

# School of Economics and Finance

## In a Small Moment: Class Size and Moral Hazard in the Mezzogiorno

Joshua D. Angrist, Erich Battistin and Daniela Vuri

Working Paper No. 747

June 2015

ISSN 1473-0278



Queen Mary  
University of London

# In a Small Moment: Class Size and Moral Hazard in the Mezzogiorno\*

Joshua D. Angrist

MIT, IZA and NBER

Erich Battistin

Queen Mary University of London, IRVAPP and IZA

Daniela Vuri

University of Rome Tor Vergata, CEIS, CESifo and IZA

June 2015

## Abstract

An instrumental variables (IV) identification strategy that exploits statutory class size caps shows significant achievement gains in smaller classes in Italian primary schools. Gains from small classes are driven mainly by schools in Southern Italy, suggesting a substantial return to class size reductions for residents of the Mezzogiorno. In addition to high unemployment and other social problems, however, the Mezzogiorno is distinguished by pervasive manipulation of standardized test scores, a finding revealed in a natural experiment that randomly assigns school monitors. IV estimates also show that small classes increase score manipulation. Dishonest scoring appears to be a consequence of teacher shirking in grade transcription, rather than cheating by either students or teachers. Estimates of a causal model for achievement with two endogenous variables, class size and score manipulation, suggest that the effects of class size on measured achievement are driven entirely by the relationship between class size and manipulation. These findings show how consequential score manipulation can arise even in assessment systems with few NCLB-style accountability concerns.

JEL Classification: C26, C31, I21, I28, J24

Keywords: test scores, education production, regression discontinuity

---

\*Special thanks go to Patrizia Falzetti, Roberto Ricci and Paolo Sestito at INVALSI for providing the achievement data used here and to INVALSI staffers Paola Giangiacomo and Valeria Tortora for advice and guidance in our work with these data. Grateful thanks also go to Gianna Barbieri, Angela Iadecola, and Daniela Di Ascenzo at the Ministry of Education (MIUR) for access to and assistance with administrative schools data. Chiara Perricone provided expert research assistance. Our thanks to David Autor, Daniele Checchi, Eric Hanushek, Andrea Ichino, Brian Jacob, Michael Lechner, Steve Machin, Derek Neal, Parag Pathak, Daniele Paserman and Jona Rockoff for helpful discussions and comments, and to seminar participants at NBER Education Fall 2013, the 2014 SOLE meeting, the University of California Irvine, Norwegian School of Economics (NHH), Padova University, IRVAPP, EUI, UCL, ISER (Essex), the CEP Labour Market Workshop, the Warwick 2014 CAGE conference, the 2014 Laax Labor Economics Workshop, the University of Rome Tor Vergata, EIEF, CEPR/IZA European Summer Symposium in Labour Economics, Tinbergen Institute and Oxford for helpful comments. This research is supported by the Einaudi Institute of Economics and Finance (EIEF) - Research Grant 2011 and by the Fondazione Bruno Kessler. Angrist thanks the Institute for Education Sciences for financial support. The views expressed here are those of the authors alone.

# 1 Introduction

School improvement efforts often focus on inputs to education production, the most important of which is staffing ratios. Parents, teachers, and policy makers look to small classes to boost learning. The question of whether changes in class size have a causal effect on achievement remains controversial, however. Regression estimates often show little gain to class size reductions, with students in larger classes sometimes appearing to do better (Hanushek, 1995). At the same time, a large randomized study, the Tennessee STAR experiment, generated evidence of substantial learning gains in smaller classes (Krueger, 1999). An investigation of longer-term effects of the STAR experiment also suggests small classes increased college attendance (Chetty et al., 2011).

Standardized tests provide the yardstick by which school quality is most often assessed and compared. As testing regimes have proliferated, however, so have concerns about the reliability and fidelity of assessment results (Neal, 2013, lays out the issues in this context). In an early empirical contribution, Jacob and Levitt (2003) documented substantial cheating in Chicago public schools, while a recent system-wide cheating scandal in Atlanta now threatens to send many school administrators to jail (Severson, 2011). Of course, students may cheat as well, especially on tests that have consequences beyond assessment. In many cases, however, the behavior of staff who administer and, in some cases, grade assessments is of primary concern. For example, Dee et al. (2011) show that scores on New York's Regents exams are often manipulated by the school staff who grade them. Score manipulation is also a potential concern in evaluations of Sweden's school choice reform (Böhlmark and Lindahl, 2013). In public school systems with weak employee performance standards, such as the Italian public school system studied here, fidelity of school staff to test administration protocols and grading standards may be especially weak.<sup>1</sup>

Our investigation of education production in Italy begins by applying the quasi-experimental research design introduced by Angrist and Lavy (1999). This design exploits variation in class size that originates in rules stipulating a class size cutoff. In Israel, with a cutoff of 40, we expect to see a single class of 40 in a grade cohort of 40, while with enrollment of 41, the

---

<sup>1</sup>De Paola et al. (2014) estimate the effects of workplace accountability on productivity in the Italian public sector. Ichino and Tabellini (2014) discuss possible benefits from organizational reform and increased choice in Italian public schools.

cohort is split into two much smaller classes. Angrist and Lavy called this “Maimonides’ Rule,” after the prominent medieval scholar Moses Maimonides, who identified a similar rule in the Talmud. Maimonides-style estimates of the effects of class size on achievement for the population of Italian second and fifth graders, most of whom attend much smaller classes than those seen in Israel, suggest a statistically significant but modest return to decreases in class size. Importantly, however, our estimated returns to class reductions are roughly three times larger in data from Southern Italy than for the rest of the country.

Italy is characterized by a sharp North-South divide along many dimensions, a fact that motivates our investigation of regional differences in class size effects. The South, known as the Mezzogiorno, is distinguished by persistently higher unemployment, lower per-capita income, higher crime rates, and lower educational attainment than are characteristic of other Italian regions.<sup>2</sup> The Mezzogiorno also lags the rest of Italy in financial development (Guiso et al., 2004), political accountability (Nannicini et al., 2013), and workplace productivity (Ichino and Maggi, 2000). Italy’s North-South divide, which is larger and more persistent than differences seen across America’s Mason-Dixon line, has been linked to cultural differences and differences in residents’ view of the role of government (Putnam et al., 1993).

Against a backdrop of relative under-development, the Mezzogiorno is also distinguished by widespread manipulation on the standardized tests given in Italian primary schools. This can be seen in Figure 1, which reproduces provincial estimates of score manipulation from the Italian Istituto Nazionale per la Valutazione del Sistema dell’Istruzione (INVALSI), a government agency charged with educational assessment. Classes in which scores are likely to have been manipulated are identified through a statistical model that looks for surprisingly high average scores, low within-class variability, and implausible missing data patterns.<sup>3</sup> Measured in this way, about 5 percent of scores are compromised, much as reported for Chicago elementary schools by Jacob and Levitt (2003). In Southern Italy, however, the proportion of compromised exams averages about 14 percent (see Table 1) and reaches 25 percent in some provinces. Further evidence on score manipulation comes from Bertoni et al.

---

<sup>2</sup>The Mezzogiorno consists of the administrative regions of Basilicata, Campania, Calabria, Puglia, Abruzzo, Molise, and the islands of Sicily and Sardinia. Italy’s 20 Administrative regions are further divided into over 100 provinces.

<sup>3</sup>The INVALSI testing program is described below and in INVALSI (2010). The INVALSI score manipulation variable identifies classes with substantially anomalous score distributions, imputing a probability of manipulation for each (see Quintano et al., 2009). Figure 1 uses this variable for the 2009-11 scores of second and fifth graders.

(2013), who analyze data generated by the random assignment of monitors sent to be present on the days tests were administered. This analysis also uncovers evidence of a substantial regional gradient in score manipulation.

In this paper, we argue that score manipulation in Italian primary schools reflects teacher behavior - specifically, dishonest transcription of students hand-written answer sheets onto machine-readable score report forms - and that teachers' time costs or disutility from dishonest score reporting increases with class size. Dishonest score reporting in this case appears to be largely a form of shirking, that is, moral hazard in grading effort, rather than cheating motivated by accountability concerns (that is, the sort of cheating discussed by Holmstrom and Milgrom, 1991). Evidence for this interpretation comes from a juxtaposition of regional patterns in the causal effects of class size on achievement and score manipulation, and from an analysis of score distributions and response patterns in classes where manipulation is likely. These patterns are estimated by exploiting the random assignment of school test monitors, who are also responsible for transcription in the schools to which they're assigned. Finally, motivated by these results, we model achievement as a function of two endogenous variables, class size and score manipulation. The model is identified by a combination of Maimonides' Rule and the monitoring policy. The resulting estimates suggest that the relationship between class size and INVALSI test scores is explained entirely by score manipulation: class size is unrelated to human capital in Italy.

Why is the fact that score manipulation explains class size effects in the Mezzogiorno of general interest? The Maimonides' Rule research design is motivated by an effort to quantify causal class size effects. This design is not guaranteed to work; Urquiola and Verhoogen (2009) show how endogenous sorting by students induces selection bias in comparisons across class size caps in Chilean private schools. By contrast, our analysis uncovers a substantive problem inherent in analyses of the causal effects of class size, regardless of research design. We show that even where the evidence that class size affects test scores is uncompromised, this need not signal increased learning in smaller classes. Our findings also provide evidence of moral hazard in a system with weak incentives. Italian teachers work in a highly regulated public sector, with virtually no risk of termination, and are subject to a pay and promotion structure that's largely independent of performance. In contrast with the unintended consequences of test-based accountability regimes, the manipulation uncovered here arises because worker

performance standards are weak. It seems fair to say that Italian moral hazard arises from a lack of accountability rather than an over-abundance of it. Finally, concerns with teacher shirking are far from unique to Italy. Clotfelter et al. (2009) discusses distributional and other consequences of American teacher absenteeism, while teacher absenteeism and other forms of shirking are a perennial concern in developing countries (see Banerjee and Duflo, 2006 and Chaudhury et al., 2006).

The rest of the paper is organized as follows. The next section presents institutional background on Italian schools and tests. Section 3 describes our data and documents the Maimonides' Rule first stage. Following a brief graphical analysis, Section 4 reports Maimonides-style estimates of effects of class size on achievement and score manipulation. Section 5 explores the nature of score manipulation by estimating score distributions and response patterns for apparent manipulators as a function of class size and item difficulty. Finally, Section 6 uses the monitoring experiment and Maimonides' Rule to jointly estimate class size and manipulation effects. This section also reviews possible threats to validity in our research design.

## 2 Background

### Italian Schools and Tests

Primary and Secondary schooling in Italy are compulsory from ages 6 to 16, with three stages: 5 years of elementary school (*scuola elementare*), lower secondary school covering grades 6-8 (*scuola media*), and high school (*scuola superiore*), which runs for 3-5 years. Schools are organized into single- or multi-unit institutions, much as a single campus might house more than one school in American public systems. Teachers are paid by seniority, without regard to qualifications, performance, or conduct.

Families apply for school admission in February, well before the beginning of the new academic year in September. Parents or legal guardians typically apply to a school in their province, located near their home. In (rare) cases of over-subscription, distance usually determines who has a first claim on seats. Rejected applicants contact other schools, mostly nearby. School principals group students into classes and assign teachers over the summer, but parents learn about class composition only in September, shortly before school starts.

At this point, parents who are unhappy with a teacher or classroom assignment are likely to find it difficult to change schools.

Italian schools have long used matriculation exams for tracking and placement in the transition from elementary to middle school and throughout high school, but standardized testing for evaluation purposes is a recent development. In 2008, INVALSI piloted voluntary assessments in elementary school; in 2009 these became compulsory for all schools and students. INVALSI assessments cover mathematics and Italian language skills in a national administration lasting two days in the Spring.<sup>4</sup> Tests are proctored by local administrators and teachers. Proctors and other teachers are expected to copy students' original responses onto machine-readable answer sheets (called *scheda risposta*), which are then sent to INVALSI. The transcription process is not entirely mechanical: some questions require teachers to interpret a student's original response as being correct, incorrect, or missing, in effect, a form of grading. Sample test items and a score sheet are included here in a brief appendix. This transcription procedure opens the door to score manipulation, as does the fact that INVALSI test administrations are typically proctored by teachers.<sup>5</sup>

## Related work

Maimonides-style empirical strategies have been used to identify class size effects in many countries, including the US (Hoxby, 2000), France (Piketty, 2004 and Gary-Bobo and Mahjoub, 2006), Norway (Bonesronning, 2003 and Leuven et al., 2008) and the Netherlands (Dobbelsteen et al., 2002). On balance these results point to modest returns to size reductions, though mostly more modest than found by Angrist and Lavy (1999) for Israel. A natural explanation for the relatively large Israeli findings is the unusually large classes characteristic of Israeli elementary schools. In line with this view, Woessmann (2005) finds a weak association between class size and achievement in a cross-country panel covering Western European school systems in which classes tend to be small.

---

<sup>4</sup>INVALSI reports school and class average scores to schools but not students. School leaders may choose to release this information to the public. Individual test scores are not reported or released. See <http://www.invalsi.it> for additional background.

<sup>5</sup>Teacher proctoring and local grading is a feature shared with other European assessments. For example, local teachers mark the UK's Key Stage 1 assessments (given in year 2, usually at age 7). Key Stage 2 assessments given at the end of elementary school (usually at age 11) are locally proctored with unannounced external monitoring and external marking (grading). See documents and links at <http://www.education.gov.uk/sta/assessment>.

The returns to class size in Italy have received little attention from researchers to date, in large part because test score data have only recently become available. Among the few Italian studies we've seen, Bratti et al. (2007) report regression estimates showing an insignificant class size effect. In an aggregate analysis, Brunello and Checchi (2005) look at the relationship between staffing ratios and educational attainment for cohorts born before 1970; they find that higher pupil-teacher ratio at the regional level are associated with higher average schooling attainment. We haven't seen other quantitative explorations of Italian class size, though Ballatore et al. (2013) use a related identification strategy to estimate the effects of the number of immigrants in the classroom on native achievement. This work doesn't consider the possible consequences of cheating or misreporting.

Jacob and Levitt (2003) and Dee et al. (2011) quantify teacher cheating on standardized assessments in Chicago and New York. The (natural) experiment that we use to identify the effects of Italian score manipulation and class size jointly was first analyzed by Bertoni et al. (2013). Our analysis of this experiment looks at monitoring effects by region, while also adjusting for features of the intervention sampling scheme not fully accounted for in earlier work. The resulting estimates suggest that the presence of external monitors sharply reduces score manipulation, and that manipulation boosts measured scores dramatically. Both of these effects are much larger in Southern Italy. Elsewhere in Italy, manipulation is relatively rare.

Scholars have documented a range of economic and behavioral differences across Italian regions. Southern Italy is characterized by low levels of social capital (Guiso et al., 2004; Guiso et al., 2010) and more widespread opportunistic or anti-social behavior (Ichino and Maggi, 2000; Ichino and Ichino, 1997). Differences along these dimensions have been used to explain persistent regional differentials in economic outcomes (Costantini and Lupi, 2006) and differences in the quality of local institutions (Putnam et al., 1993). Finally, as noted in the introduction, our work connects with research on teacher shirking around the world.



### 3 Data and First Stage

#### Data and descriptive statistics

The standardized test score data used in this study come from INVALSI's testing program in Italian elementary schools in the 2009/10, 2010/11, and 2011/12 school years. Raw scores indicate the number of correct answers; for the purposes of regression and two-stage least squares (2SLS) estimation, we standardized these by subject, year of survey, and grade to have zero mean and unit variance. Data on test scores were matched to administrative and survey information describing institutions, schools, classes, and students. Class size can be measured by administrative enrollment counts at the beginning of the school year as well as the number of test-takers (we use the former). Student data include gender, citizenship, and information on parents' employment status and educational background. These data are collected as part of test administration and supposed to be provided by school staff when scores are submitted. Fewer than 10 percent of Italian primary and secondary school students who attend private schools are omitted from this study.

Our statistical analysis focuses on class-level averages since this is the aggregation level at which the regressor of interest varies. The empirical analysis is restricted to classes with more than the minimum number of students set by law (10 before 2010 and 15 from 2011). This selection rule eliminates classes in the least populated areas of the country, mostly mountainous areas and small islands. We also drop schools enrolling more than 160 students in a grade, as these are above the threshold where Maimonides' Rule is likely to matter (this size cutoff trims classes above the 99th percentile of the enrollment-weighted class size distribution).

The resulting matched file includes about 70,000 classes in each of the two grades covered by our three-year window (these are repeated cross-sections; the data structure doesn't follow the same classes over time). Table 1 shows descriptive statistics for the estimation sample, separately by grade. Statistics are reported at the class level in Panel A, at the school level in Panel B, and at the institution level in Panel C. Class size averages around 20 in both grades, and is slightly lower in the South. The score means reported in Panel A give the class average percent correct. Scores are higher in language than in math and higher in grade 5 than in grade 2. The table also shows averages for an indicator of score manipulation variable

that we've constructed (Section 3, below, explains how this was done). Manipulation rates are higher in the South and in math.

### Maimonides in Italy

Our identification strategy exploits minimum and maximum class sizes for Italy (these rules are a consequence of a regulation known as *Decreto Ministeriale 331/98*). Until the 2008/09 school year, primary school classes were subject to a minimum size of 10 and capped at 25. Grade enrollment beyond 25 or a multiple thereof usually prompted the addition of a class. The rule allows exceptions, however. Principals can reduce the size of classes attended by one or more disabled students, and schools in mountainous or remote areas are allowed to open classes with fewer than 10 students. Finally, the law allows a 10% deviations from the maximum in either direction (that is, the Ministry of Education will usually fund an additional class when enrollment exceeds 22 and typically requires a new class when average enrollment would otherwise exceed 28). A 2009/10 reform increased the nominal maximum to 27, with a minimum size of 15, again with a tolerance of 10% (promulgated through *Decreto del Presidente della Repubblica 81/2009*). This reform was rolled out one grade per year, starting with grade 1. In our data, second graders in 2009/10 and fifth graders in any year are subject to the old rule, while second graders in 2010/11 and 2011/12 are subject to the new rule.

Ignoring discretionary deviations near cutoffs, Maimonides' Rule predicts class size to be a non-linear and discontinuous function of enrollment. Writing  $f_{igkt}$  for the predicted size of class  $i$  in grade  $g$  at school  $k$  in year  $t$ , we have:

$$f_{igkt} = \frac{r_{gkt}}{[\text{int}((r_{gkt} - 1) / c_{gt}) + 1]}, \quad (1)$$

where  $r_{gkt}$  is beginning-of-the-year grade enrollment at school  $k$ ,  $c_{gt}$  is the cap in effect that year (25 or 27) in grade  $g$ , and  $\text{int}(x)$  is the largest integer smaller than or equal to  $x$ . Figure 2 and Figure 3 plot average class size and  $f_{igkt}$  against enrollment in grade, separately for pre- and post-reform periods. Plotted points show the average actual class size at each value of enrollment. Actual class size follows predicted class size reasonably closely for enrollments below about 75, especially in the pre-reform period. Theoretical sharp corners in the class size/enrollment relationship are rounded by the soft nature of the rule. Many classes are

split before reaching the theoretical maximum of 25. Earlier-than-mandated splits occur more often as enrollment increases. In the post-reform period, class size tracks the rule generated by the new cap of 27 poorly once enrollment exceeds about 70.

## Measuring Manipulation

Our score manipulation variable is a function of extreme values, the within-class average and standard deviation of test scores, the number of missing items, and a Herfindahl index of the share of students with similar response patterns. These indicators are used as inputs for a cluster analysis that flags as suspicious classes with abnormally high performance, an unusually small dispersion of scores, an unusually low proportion of missing items, and high concentration in response patterns. This procedure yields class-level indicators of compromised scores, separately for math and language. The resulting manipulation indicator is similar to the manipulation variable used by Quintano et al. (2009) and in INVALSI publications (e.g., INVALSI, 2010). The INVALSI version generates a continuous class-level probability of manipulation. The procedure used here generates a dummy variable indicating classes where score manipulation seems likely. Methods and formulas used to classify score manipulation are detailed further in the appendix.<sup>6</sup>

## 4 Class Size Effects: Achievement and Manipulation

### Graphical Analysis

We begin with plots that capture class size effects near enrollment cutoffs. The first in this sequence, Figure 4, documents the relationship between cutoffs (multiples of 25 or 27) and class size. This figure was constructed from a sample of classes at schools with enrollment that falls in a  $[-12,12]$  window around the first four cutoffs shown in Figure 2 and Figure 3. Enrollment values in each window are centered to be zero at the relevant cutoff. The y-axis shows average class size conditional on the centered enrollment value shown on the x-axis,

---

<sup>6</sup>Our procedure also follows Jacob and Levitt (2003) in inferring score manipulation from patterns of answers within and across tests in a classroom. Jacob and Levitt (2003) also compare test scores over time, looking for anomalous changes. Values in the upper tail of the Jacob-Levitt suspicious answer index are highly predictive of their cheating variable in the cross section. Our main results are unchanged when manipulation is measured continuously. A binary indicator leads to parsimonious models and easily interpreted estimates, however, while also facilitating the discussion of misclassification bias in section 6.2.

reported as a 3-point moving average. Figure 4 also plots fitted values generated by local linear regressions (LLR) fits to class-level data. In this context, the LLR smoother uses data on one side of the cutoff only, smoothed with an edge kernel and Imbens and Kalyanaraman (2012) bandwidth.<sup>7</sup>

In view of the 2-3 student tolerance around the cutoff for the addition of a class, enrollment within two points of the cutoff is excluded from the local linear fit. As a result of this tolerance, class size can be expected to decline at enrollment values shortly before the cutoff and to continue to decline thereafter. Consistent with this expectation, the figure shows a clear drop at the cutoff, with the sharpness of the break moderated by values near the cutoff. Class size is minimized at about 3-5 students to the right of the cutoff instead of immediately after, as we would expect were Maimonides Rule to be tightly enforced. The parametric identification strategy detailed below exploits both the discontinuous variation in class size generated when enrollment moves across cutoffs and changes in slope as a cohort is divided into classes more finely. Looking only at points immediately adjacent to the cutoff, the change in size generating by moving across a cutoff is on the order of 2-3 students.

In data from the South, math and language scores plotted as a function of enrollment values near Maimonides cutoffs show a jump that mirrors the drop in size seen at Maimonides cutoffs, but there is little evidence of such a jump in schools outside the South. This pattern is documented in Figure 5, which plots math and language scores against enrollment in a format paralleling that of Figure 4. The reduced-form achievement drop for schools in Southern Italy is about 0.02 standard deviations (hereafter,  $\sigma$ ). Assuming this reduced-form change in test scores in the neighborhood of Maimonides cutoffs is driven by a causal class size effect, the implied return to a one-student reduction in class size is about  $0.01\sigma$  in Southern Italy (this comes from dividing 0.02 by a rough first stage of about 2). The absence of a jump in scores at cutoffs in data from schools elsewhere in the country suggests that outside the South class size reductions leave scores unchanged.

Score manipulation also varies as a function of enrollment in the neighborhood of class size cutoffs, with a pattern much like that seen for achievement. This is apparent in Figure 6, which puts the proportion of classes identified as having compromised scores on the y-axis, in

---

<sup>7</sup>The figures here plot residuals from a regression of class size on the controls included in equation (2), below.

a format like that used for Figure 4 and Figure 5. Mirroring the pattern of achievement effects, a discontinuity in score manipulation rates emerges most clearly for schools in Southern Italy. This pattern suggests that the achievement gains generated by class size in Figure 5 may reflect the manipulation behavior captured in Figure 6. A possible caveat here is the role of mismeasured manipulation might have in generating this pattern. We explore implications of misclassification for our 2SLS estimates in a separate section, following these estimates. We note here, however, that classification error is unlikely to change discontinuously at class size cutoffs. Moreover, the fact that manipulation is essentially smooth through the cutoff for schools outside the South weighs against purely mechanical explanations of the pattern in Figure 6.

### Empirical Framework

Figure 5 suggests that variation in class size near Maimonides cutoffs can be used to identify class size effects in a non-parametric fuzzy regression discontinuity (RD) framework. In what follows, however, we opt for parametric models that exploit variation in enrollment due to changes in the slope of the relationship between enrollment and class size, as well as discontinuities. The parametric strategy gains power by combining features of both RD and regression kink designs, while easily accommodating models with multiple endogenous variables and covariates.<sup>8</sup>

Our parametric framework models  $y_{igkt}$ , the average outcome score in class  $i$  in grade  $g$  at school  $k$  in year  $t$ , as a polynomial function of the running variable,  $r_{gkt}$ , and class size,  $s_{igkt}$ . With quadratic running variable controls, the specification pooling grades and years can be written:

$$y_{igkt} = \rho_0(t, g) + \beta s_{igkt} + \rho_1 r_{gkt} + \rho_2 r_{gkt}^2 + \epsilon_{igkt}, \quad (2)$$

where  $\rho_0(t, g)$  is shorthand for a full set of year and grade effects. This model also controls for the demographic variables described in Table 1, as well as the stratification variables used in the monitoring experiment to increase precision in the estimates.<sup>9</sup> Standard errors are clustered by institution, which we reckon to be a conservative strategy in this context.

---

<sup>8</sup>Card et al. (2012) discuss nonparametric identification in the regression kink design.

<sup>9</sup>Control variables include proportion female in the class, the proportion of immigrants, the proportion of students whose father is a high school graduate, have unemployed mothers, have mothers not in the labor force, have employed mothers, and dummies for missing values for these variables. Stratification controls consist of total enrollment in grade, region dummies, and the interaction between enrollment and region.

The instrument used for 2SLS estimation of equation (2) is  $f_{igkt}$ , as defined in equation (1). To document the sensitivity of findings to specification of running variable controls, we also report results from models that include a full set of cutoff-segment (window) main effects, while allowing the quadratic control function to differ across segments. We refer to this as the interacted specification.<sup>10</sup> The corresponding OLS estimates for models without interacted running variable controls are shown as a benchmark.

### Parametric Estimates of Class Size Effects

OLS estimates of equation (2) show a negative correlation between class size and achievement for schools in the Northern and Central regions, but not in the South (class size effects are scaled for a 10 student change). Larger classes are associated with somewhat higher language scores in the South while Southern class sizes appear to be unrelated to achievement in math. These estimates can be seen in columns 1-3 of Table 2.

2SLS estimates using Maimonides' Rule, reported in columns 4-9 of Table 2, suggest that larger classes reduce achievement in both math and language. The associated first stage estimates, which can be seen in Appendix Table A1, show that predicted class size increases actual class size with a coefficient around one-half when regions are pooled, with a first stage effect of 0.43 in the South and 0.55 elsewhere. 2SLS estimates for Southern schools, implying something on the order of a  $0.10\sigma$  achievement gain for a 10-student reduction, are 2-3 times larger than the corresponding estimates for schools outside the South. The 2SLS estimates are reasonably precise; only estimates of the interacted specification for language scores from non-Southern schools fall short of conventional levels of statistical significance. On balance, the results in Table 2 indicate a substantial achievement payoff to class size reductions, though the gains here are not as large as those reported by Angrist and Lavy (1999) for Israel. A substantive explanation for this difference in findings might be concavity in the relationship between class size and achievement, combined with Italy's much smaller average class sizes.

The estimates in Table 3 suggest that the causal effect of class size on measured achieve-

---

<sup>10</sup>Pre-reform segments cover the intervals 10-37, 38-62, 63-87, 88-112, 113-137, and 138-159; post-reform segments cover the intervals 15-40, 41-67, 68-94, 95-121, and 122-159. These segments cover intervals of width +/- 12 in the pre-reform period and +/-13 in the post-reform period, with modifications at the lower and upper segments to include a few larger and smaller values.

ment reported in Table 2 need not reflect more learning in smaller classes. This table reports estimates from specifications identical to those used to construct the estimates in Table 2, with the modification that a class-level score manipulation indicator replaces achievement as an outcome. The 2SLS estimates in columns 4-9 show a large and precisely-estimated negative effect of class size on manipulation rates, with effects on the order of 4-6 percentage points for a 10-student class size increase in the South. Estimates for schools outside the South also show a negative relationship between class size and score manipulation, though here the estimated effects are much smaller and significantly different from zero in only one case (language scores from the non-interacted specification). Interestingly, OLS estimates of effect of class size on score manipulation, though smaller in magnitude, reflect the same negative effects as 2SLS. This suggests that the relationship between class size and manipulation may have a mechanical component, less affected by the sort of selection bias that affects OLS estimates of the corresponding achievement relation.

## 5 The Nature of Manipulation

### 5.1 The Monitoring Experiment

We turn here to the questions of who manipulates scores and why manipulators act as they do. These questions are explored with the help of INVALSI's monitoring intervention, which provides an independent source of quasi-experimental variation in score manipulation, unrelated to Maimonides' Rule. In an effort to increase test reliability, INVALSI randomly selects institutions to be observed by an external monitor on test days. Institutions are sampled for monitoring with a probability proportional to grade enrollment in the year of the test. Sampling is also stratified by regions. One class is selected for monitoring in sampled institutions with grade enrollment below 100. Two classes are selected in remaining institutions (randomness of within-institution monitoring appears to have been compromised in practice). Regional education offices select monitors from a pool of mainly retired teachers and principals who've not worked in the towns or at the schools they are assigned to monitor for at least two years. Monitors supervise test administration and encourage compliance with INVALSI testing standards. Monitors also supervise score sheet transcription onto the *scheda*

*risposta*, a burdensome clerical task that’s meant to be completed by the end of the test day. Tests without monitors are proctored by local school staff (though the math teacher for a given grade is not supposed to be assigned to proctor that grade’s test and so on). Although monitors are salaried, teachers are required to grade exams and transcribe scores outside school hours, tasks for which they receive no additional pay. This extra uncompensated work opens the door to moral hazard in teacher behavior (the monitoring policy was introduced to counteract this).<sup>11</sup>

Table 4 documents balance across institutions with and without randomly assigned monitors. Specifically, this table shows regression-adjusted treatment-control differences from models that control for strata in the monitoring sample design. These specifications include a full set of region dummies and a linear function of institutional grade enrollment that varies by regions (standard errors are clustered by institution). Administrative variables – generated as a by-product of school administration and INVALSI testing – are well-balanced across groups, as can be seen in the small and insignificant coefficient estimates reported in Panel A of the table.<sup>12</sup>

Demographic data and other information provided by school staff, such as parental information, show evidence of imbalance. This seems likely to reflect the influence of monitoring on data quality, rather than a problem with the experimental design or implementation. The hypothesis that monitors induced more careful data reporting by staff is supported by the large treatment-control differential in missing data rates documented at the bottom of the table. Among other salutary effects, randomly assigned monitors reduce item non-response by as much as three percentage points, as can be seen in Panel C of Table 4. Monitoring effects on data quality at class size cutoffs are discussed in Section 6.2.

The presence of institutional monitors reduces score manipulation considerably. This is apparent in the estimated monitoring effects shown in columns 1-3 of Table 5, which reports estimates showing that monitoring reduces manipulation rates by about 3 percentage points for Italy, with effects twice as large in the South. These estimates come from models

---

<sup>11</sup>Econometric estimates of the effect of monitoring on score manipulation were first reported by INVALSI (2010). We replicate some of these earlier findings, as well as those reported by Bertoni et al. (2013).

<sup>12</sup>Bertoni et al. (2013) mistakenly treated institutions as schools. Their identification strategy also presumes random assignment of classroom monitors within institutions, but we find that monitors are much more likely to be assigned to large classes, probably a consequence of that fact that in most institutions only one class is monitored.



similar to those used to check covariate balance with a score manipulation indicator replacing covariates as the dependent variable. Monitoring also reduces language scores by  $0.08\sigma$ , while the estimated monitoring effect on math scores is about  $-0.11\sigma$ . Effects of monitoring in the South range from  $-0.13\sigma$  for language to  $-0.18\sigma$  for math, estimates that appear in column 6 of the table.

The estimates reported in columns 1-3 and 4-6 of Table 5 constitute the first stage and reduced form for a model that uses the assignment of monitors as an instrument for the effects of score manipulation on test scores. Dividing reduced form estimates by the corresponding first stage estimates produces second stage manipulation effects of about  $3\sigma$  for the South, with even larger second stage estimates for the North. These effects seem implausibly large, implying a boost in scores that exceeds the range of the dependent variable in some cases. It also seems likely, however, that the score manipulation variable used to construct the corresponding first stage effects is an imperfect measure of actual manipulation behavior. Because classification error attenuates first stage estimates in this context, the resulting second stage estimates are proportionally inflated. This and other implications of missclassification are discussed in Section 6, below, after reviewing estimation results that simultaneously capture class size and manipulation effects on test scores.<sup>13</sup>

The central role of monitoring in the link between class size and measured achievement link is supported by Table 6, which reports 2SLS estimates of class size effects on test scores, separately for institutions with and without INVALSI monitors. These estimates reveal a strong negative effect of class size on achievement - but only in the absence of monitoring. These findings suggest that in the absence of monitors, larger class sizes facilitate or encourage manipulation.

## 5.2 The Anatomy of Manipulation

The facts that monitoring reduces score manipulation and that score manipulation *decreases* with class size strongly suggests that teachers are the source of manipulation and not students. Honest teacher-proctors should have the same deterrent effect as external monitors on cheating students: both are likely to catch cheaters, perhaps teachers even more so if they

---

<sup>13</sup>Using survey data on exam day experiences and perceptions, Bertoni et al. (2013) find no direct effects of monitors on fifth graders' feelings or motivation.

recognize cheating more readily. Moreover, any class size effect on student cheating is likely to be positive, that is, larger classes should facilitate student cheating by making cheating harder to detect. Results in Table 3 showing that score manipulation decreases with class size therefore weigh against student cheating as well. Finally, because individual test scores are never disclosed, it's hard to see why students might care to cheat. At the same time, the fact that teachers are required to transcribe scores – except when monitors do it for them – provides a natural opportunity for score misreporting.

What might motivate score manipulation by teachers? Cheating on standardized assessments often appears to reflect a desire to boost scores in the face of school accountability considerations, as has been seen in a number of US school districts. Jacob and Levitt (2003), for example, document a pattern of accountability-motivated score manipulation in the Chicago Public Schools. Yet, school accountability in Italy is weak, with no systematic review of school-level or teacher-related achievement outcomes by principals, parents, or government officials.

Although we have no direct evidence on teacher motivation, a comparison of the pattern of item-level scores in classrooms where score manipulation is likely to the pattern in classrooms where manipulation is unlikely provides evidence on the question. Other things equal, manipulators motivated by accountability considerations would probably prefer not to get caught. Such manipulation therefore seems likely to take the form of supplying answers or hints on difficult items while proctoring exams. Specifically, a desire to limit exposure and cheat efficiently should induce cheating that focuses on items where scores are otherwise likely to be low. In other words, accountability considerations seem likely to make manipulation behavior that targets difficult items especially attractive.

To investigate the relationship between manipulation and item difficulty, we measure difficulty using the percent correct on item  $j$ , denoted by  $p_j$ , for monitored institutions in Veneto, a province where manipulation rates are very low. When accountability concerns are paramount, items with high  $p_j$  should be graded honestly by manipulators, leaving a strong positive relationship between outcome scores and difficulty when  $p_j$  is high. At the same time, as  $p_j$  falls, selective manipulation should flatten the score gradient, making the overall relationship between manipulated scores and  $p_j$  convex. By contrast, mechanical manipulation behavior - such as copying entire answer sheets - should push scores on all

items up to the same high level.

Because monitors are randomly assigned to institutions and not classes, we identify latent score distributions in classes with and without manipulation using the instrumental variables methods developed by Abadie (2002) (an application of this approach to school reform treatment effects appears in Angrist et al. 2013). Specifically, each classroom is taken to have two potential score distributions for each item,  $j$ , one revealed in the presence of manipulation ( $y_{igkt}^j(1)$ ) and one revealed otherwise ( $y_{igkt}^j(0)$ ). Observed scores in class  $i$  on item  $j$ , denoted by  $y_{igkt}^j$ , are determined by

$$y_{igkt}^j = (1 - m_{igkt})y_{igkt}^j(0) + m_{igkt}y_{igkt}^j(1),$$

where  $m_{igkt}$  is a class-level manipulation indicator (there are about 45 items per year/grade/subject).

Monitoring,  $M_{igkt}$ , provides an exogenous shock to manipulation behavior. The difficulty gradient for manipulated and non-manipulated scores is identified by 2SLS estimation of the parameters  $\beta_1^j$  and  $\beta_0^j$  in models of the form

$$\begin{aligned} y_{igkt}^j m_{igkt} &= \rho_1(t, g) + \beta_1^j m_{igkt} + \epsilon_{igkt}, \\ y_{igkt}^j (1 - m_{igkt}) &= \rho_0(t, g) + \beta_0^j (1 - m_{igkt}) + \epsilon_{igkt}, \end{aligned}$$

for item, using data from the South, where manipulation is prevalent. The class-level manipulation indicators,  $m_{igkt}$  and  $1 - m_{igkt}$ , are treated as endogenous and instrumented by randomly assigned institutional monitoring,  $M_{igkt}$ . The resulting estimates of  $\beta_1^j$  capture potential scores on item  $j$  under manipulation for complying classes, that is, for classes in which we can expect manipulation in the absence of monitoring and honest scoring otherwise. Similarly, the parameter  $\beta_0^j$  is the average potential score on item  $j$  without manipulation for the same classes.

Manipulation changes the relationship between item difficulty and test scores markedly, pushing an otherwise steep difficulty gradient up to a high level, with scores uniformly close to 100 percent correct. This can be seen in Figure 7, which plots our 2SLS estimates of  $\beta_1^j$  and  $\beta_0^j$  against  $p_j$ . This figure also shows a least squares fit to the relationship between scores and difficulty, weighted by the precision of the item-level estimates. Manipulation for accountability purposes seems unlikely to produce the essentially linear and nearly flat relationship between item difficulty and manipulators' test scores apparent in Figure 7.

Figure 7 also distinguishes items by the level of effort required for transcription. Some items are transcribed quickly and easily onto the machine-readable *scheda risposta*, but others require thought and judgement; transcription of these items is more of a grading exercise than a copying task. Examples of high-grading-effort items are given in the appendix. In view of this difference in effort, teachers might target high-grading-effort items for manipulation, answering randomly or copying correct answers to such items. If manipulators focus on high-grading-effort items, we should see large score differences by manipulation status for such items only. A comparison of the left and right panels in Figure 7, however, offers little evidence of such targeted manipulation behavior: conditional on difficulty, the difference in scores between manipulators and non-manipulators is similar for high- and low-grading-effort items.

The item-level analysis in Figure 7 offers little evidence of selective manipulation based either on item difficulty or grading effort. The fact that manipulated scores are well above honest scores also makes pervasive random transcription unlikely. What sort of behavior is consistent with the patterns apparent in the figure? In this case, the simplest story seems most likely: manipulating teachers would appear to forgo honest transcription entirely, copying entire answer sheets, without regard to item characteristics (a form of dishonest reporting akin to “curbstoning” in survey research). Wholesale curbstoning, a strategy that minimizes transcription or grading effort while maintaining high levels of achievement, seems to be the primary force behind score manipulation.

Finally, we ask how score manipulation through curbstoning might be linked with class size. Here too, a simple explanation seems plausible and is supported by estimates of class size effects on score distribution. Specifically, the estimates plotted in Figure 8 show 2SLS estimates of the effects of class size on a family of indicators for test scores that fall in 5-point bands (models here are the same as those used to construct the 2SLS estimates reported in Table 2). The figure reveals a clear large-class effect in the form of a shift from very high scores to low scores, with little change through the middle of the score distribution. This pattern supports our conjecture that large classes reduce scores by deterring or reducing curbstoning and, perhaps, by making curbstoning less accurate when it occurs. The forces behind this pattern include the fact that the number of teachers proctoring and transcribing exams probably increases in larger classes, limiting manipulation through peer monitoring.

At the same time, curbstoning of correct answers is probably less accurate in large classes. Of course, transcription accuracy may fall with class size without regard to cheating. Weighting against a pure accuracy-in-transcription effect, however, is the fact that the relationship between class size and test scores disappears once manipulation is taken into account. Honest transcribers would appear to do this work accurately, while shirkers may grow careless as the transcription workload grows.

## 6 Manipulation Explains Class Size Effects

### 6.1 Estimates with Two Endogenous Variables

The estimates in Table 2, Table 3, and Table 5 motivate a causal model in which achievement depends on class size ( $s_{igkt}$ ) and score manipulation ( $m_{igkt}$ ), both treated as endogenous variables to be instrumented. This model can be written:

$$y_{igkt} = \rho_0(t, g) + \beta_1 s_{igkt} + \beta_2 m_{igkt} + \rho_1 r_{gkt} + \rho_2 r_{gkt}^2 + \eta_{igkt}, \quad (3)$$

where  $\rho_0(t, g)$  is again a shorthand for year and grade effects. We interpret equation (3) as describing the average achievement that would be revealed by alternative assignments of class size,  $s_{igkt}$ , in an experiment that holds  $m_{igkt}$  fixed. This model likewise describes causal effects of changing score manipulation rates in an experiment that holds class size fixed. In other words, (3) is a model for potential outcomes indexed against two jointly manipulable treatments.

We estimate equation (3) by 2SLS in a setup that includes the same covariates that appear in the models used to construct the estimates reported in Table 2. The instrument list contains Maimonides' Rule ( $f_{igkt}$ ) and a dummy indicating classes at institutions with randomly assigned monitors,  $M_{igkt}$ . The first-stage equations associated with these two instruments can be written:

$$s_{igkt} = \lambda_{10}(t, g) + \mu_{11} f_{igkt} + \mu_{12} M_{igkt} + \lambda_{11} r_{gkt} + \lambda_{12} r_{gkt}^2 + \xi_{ik}, \quad (4)$$

$$m_{igkt} = \lambda_{20}(t, g) + \mu_{21} f_{igkt} + \mu_{22} M_{igkt} + \lambda_{21} r_{gkt} + \lambda_{22} r_{gkt}^2 + \upsilon_{ik}, \quad (5)$$

where  $\lambda_{10}(t, g)$  and  $\lambda_{20}(t, g)$  are shorthand for first-stage year and grade effects. First stage estimates, reported in Table 7, show both a monitoring and a Maimonides' Rule effect on

score manipulation, both of which are considerably more pronounced in the South. The Maimonides first stage for class size remains at around one-half, while the presence of a monitor at institution is unrelated to class size. This is consistent with the hypothesis that monitors are randomly assigned to institutions.

The 2SLS estimates of  $\beta_2$  in equation (3), reported in Table 8, show large effects of manipulation on test scores. At the same time, this table reports small and mostly insignificant estimates of  $\beta_1$ , the coefficient on class size in the multivariate model. In an effort to boost the precision of these estimates, we estimate over-identified models that add four dummies for values of the running variable that fall within 10% of each cutoff, a specification motivated by the non-parametric first stage captured in Figure 4.<sup>14</sup> The most precise of the estimated zeros reported in Table 8, generated by the over-identified specification for Italy as a whole, run no larger than 0.022, with an estimated standard error of 0.015 (for a 10 student increase in class size); these appear in column 4. It's also worth noting that the over-identification p-values associated with these estimates are far from conventional significance levels.

We also report 2SLS estimates adding an interaction term,  $s_{igkt} * m_{igkt}$ , to equation (3), using  $f_{igkt} * M_{igkt}$  and the extra dummy instruments interacted with  $M_{igkt}$  as excluded instruments. This specification is motivated by the idea that class size may matter only in a low-manipulation subsample, while an additive model like equation (3) may miss this. There is little evidence for interactions, however: the estimated interaction effects, reported in columns 7-9 of Table 8 are not significantly different from zero.

The most important findings in Table 8 are the small and insignificant class size effects for the Mezzogiorno, a result that contrasts with the much larger and statistically significant class size effects for the same area reported in Table 2. In column 9 of the latter table, for example, a 10 student reduction in class size is estimated to boost achievement by  $0.10\sigma$  or more. The corresponding multivariate estimates in Table 8 are of the opposite sign, showing that larger classes increase achievement, though not by very much. The over-identified estimates come with estimated standard errors ranging from about 0.02 to 0.04, so that the estimated class size effects in Table 2 fall well outside the estimated confidence intervals associated with the multivariate estimates. It seems reasonable, therefore, to interpret the estimated class effects in Table 8 as precise zeros. This in turn aligns with an interpretation of the return to class

---

<sup>14</sup>First stage estimates for the over-identified model appear in Appendix Table A2.

size in Italy as due entirely to the causal effect of class size on score manipulation, most likely by teachers.

## 6.2 Threats to validity

We briefly consider three possible threats to validity in our research design. An initial concern comes from the fact that one of the four indicators used to construct the score manipulation dummy, unusually high average scores, may be connected to the outcome of interest for reasons unrelated to manipulation. We therefore constructed a manipulation variable excluding this component. RD estimates of the relationship between class size, score manipulation, and achievement, are largely unaffected by this change. Two other concerns relate to measurement error in score manipulation and potentially endogenous sorting around class size cutoffs.

### Score manipulation with misclassification

The large 2SLS estimates of manipulation effects in Table 8 reflect attenuation bias in first stage estimates if score manipulation is misreported. We show here that as long as misclassification rates are independent of the instruments, mismeasurement of manipulation leaves 2SLS estimates of *class size effects* in the multivariate model unaffected. We show this in the context of a simplified version of the multivariate model, which can be written with a class subscript as:

$$y_i = \rho_0 + \beta_1 s_i + \beta_2 m_i^* + \zeta_i, \quad (6)$$

where instruments are assumed to be uncorrelated with the error,  $\zeta_i$ , as in equation (3). Here,  $m_i^*$  is an accurate score manipulation dummy for class  $i$ , while  $m_i$  is observed score manipulation as before.

Let  $z_i = [f_i \ M_i]'$  denote the vector of instruments. Assuming that classification rates are independent of the instruments conditional on  $m_i^*$ , we can write:

$$m_i = (1 - \pi_0) + (\pi_0 + \pi_1 - 1)m_i^* + \omega_i, \quad (7)$$

where the residual,  $\omega_i$ , is defined by:

$$\omega_i = m_i - E[m_i | z_i, m_i^*],$$

and  $\pi_d$ , the probability that score manipulation is correctly detected, satisfies:

$$P[m_i = d|z_i, m_i^* = d] = P[m_i = d|m_i^* = d] = \pi_d, \quad (8)$$

for  $d = 0, 1$ . Note that  $E[z_i\omega_i] = 0$  by definition of  $\omega_i$ .

Using (7) to substitute for  $m_i^*$ , equation (6) can be rewritten:

$$y_i = \left[ \rho_0 - \frac{\beta_2(1 - \pi_0)}{\pi_0 + \pi_1 - 1} \right] + \beta_1 s_i + \left[ \frac{\beta_2}{\pi_0 + \pi_1 - 1} \right] m_i + \left[ \zeta_i - \beta_2 \frac{\omega_i}{\pi_0 + \pi_1 - 1} \right]. \quad (9)$$

We assume that the  $\pi_d$ 's are strictly greater than 0.5, so that reported score manipulation is a better indicator of actual manipulation than a coin toss. This ensures that the coefficient on  $m_i$  in (9) is finite and has the same sign as  $\beta_2$ . The 2SLS estimate of the coefficient on reported score manipulation is therefore biased upward, since  $\pi_0 + \pi_1 - 1$  is strictly between 0 and 1 given these assumptions. This implies that estimates of  $\beta_2$  for the North/Centre region (columns 2, 5 and 8 of Table 8), where score manipulation is lower and therefore misclassification is higher, are more inflated than in the South. Most importantly, because the feasible estimating equation (9) has a residual uncorrelated with the instruments and the coefficient on class size is unchanged in this model, misclassification of the sort described by (8) leaves estimates of the class size coefficient,  $\beta_1$ , unchanged. Similar results for the consequences of classification error under the same assumptions appear in Kane et al. (1999), Mahajan (2006), and Lewbel (2007), among others, though our work focuses on the consequences for the coefficient on a variable subject to error rather than implications for other regressors in the model.<sup>15</sup>

### Sorting near cutoffs

The Maimonides research design identifies causal class size effects assuming that, after adjusting for secular effects of the running variable, predicted class size ( $f_{igkt}$ ) is unrelated to student or school characteristics. As in other RD-type designs, sorting around cutoffs poses a potential threat to this assumption. Urquiola and Verhoogen (2009) and Baker and Paser-

---

<sup>15</sup>We can learn whether 2SLS estimates of the coefficient on  $m_i$ , that is, the size of the estimated manipulation effects, are plausible by experimenting with data from an area where manipulation rates are low and assuming that true manipulators earn perfect scores. We use data from Veneto, the region with the lowest score manipulation rate in Italy, to estimate  $\beta_2$  in this scenario by picking 20% of classes at random and recoding scores for this group to be 100. The resulting estimates of  $\beta_2$  come out at around  $2.25\sigma$ . Taking this as a benchmark, the manipulation effects in Table 8 are consistent with values of  $\pi_j$  around .8 for Italy (since  $\frac{2.25}{2 \times .8 - 1} = 3.75$ ), though the implied  $\pi_j$ 's are closer to .65 for math scores outside the South. These rates seem like reasonable descriptions of the classification process.



man (2013) note that discontinuities in student characteristics near Maimonides cutoffs can arise if parents or school authorities try to shift enrollment to schools where expected class size is small. In our setting, however, an evaluation of the sorting hypothesis is complicated by the link between Maimonides' Rule and score manipulation documented in Table 7. The fact that Maimonides' Rule predicts score manipulation, especially in the South, generates the results in Table 8. An important channel for the link between Maimonides' Rule and manipulation is the fact that monitoring rates are lower in small classes. If the behavior driving manipulation also affects data quality, a conjecture supported by the effects of monitoring on data quality seen in Table 4, we might expect Maimonides' Rule to be related to covariates for the same reason that monitoring is related to covariates.

This expectation is borne out by Table 9, which reports estimates of the link between Maimonides' Rule and covariates in a format paralleling that of Table 4. These estimates come from the reduced form specifications used to generate the 2SLS estimates reported in Table 2, after replacing scores with covariates on the left hand side. The pattern of covariate imbalance in Table 9 mirrors that in Table 4: covariates affected by monitoring are also correlated with Maimonides' Rule, while administrative variables that are unrelated to monitoring are largely orthogonal to Maimonides' Rule. Tables 4 and 9 also reflect similar regional differences in the degree of covariate imbalance, with considerably more imbalance in the South. Additional evidence suggesting that the link between covariates is a data quality effect unrelated to sorting appears in Appendix Table A3. This table shows that the  $f_{igt}$  is largely unrelated to covariates in schools with monitors, where manipulation is considerably diminished (though not necessarily eliminated, since some classes in monitored institutions remain unmonitored).

## 7 Summary and Directions for Further Work

The causal effects of class size on Italian primary schoolers' test scores are identified by quasi-experimental variation arising from Italy's version of Maimonides' Rule. The resulting estimates show small classes boost test scores in Southern provinces, an area known as the Mezzogiorno, but not elsewhere. Analyses of data on score manipulation and a randomized institution monitoring experiment reveal substantial manipulation in the Mezzogiorno, most

likely by teachers. For a variety of institutional and behavioral reasons, teacher score manipulation is inhibited by larger classes as well as by external monitoring. Estimates of a model that jointly captures the causal effects of class size and score manipulation on measured achievement suggest the returns to class size in the Mezzogiorno are explained by the causal effects of class size on score manipulation, with no apparent gains in learning. These findings show how class size effects can be misleading even where internal validity is probably not an issue. Our results also show how score manipulation can arise as a result of shirking in an institutional setting where standardized assessments are largely divorced from accountability.

These findings raise a number of questions, including those of why teacher manipulation is so much more prevalent in the Mezzogiorno, and what can be done to enhance accurate assessment in Italy and elsewhere. Manipulation in the Mezzogiorno arises in part from local exam proctoring and local transcription of answer sheets, a strategy meant to lower costs. New York's venerable Regent's exams were also graded locally until 2013, an arrangement that likewise appears to have facilitated score manipulation. Moreover, as with INVALSI assessments, manipulation of Regent's scores appears to be unrelated to NCLB-style accountability pressure (Dee et al., 2011). By contrast, the UK's Key Stage 2 primary-level assessments are marked by external examiners, a costly but probably worthwhile effort.<sup>16</sup> It's also worth asking why class size reductions fail to enhance learning in Italy, while evidence from the US, Israel, and a number of other countries suggest class size reductions often increase learning. We hope to address these questions in future work.

---

<sup>16</sup>See [https://home.edexcelgateway.com/pages/job\\_search\\_view.aspx?jobId=537](https://home.edexcelgateway.com/pages/job_search_view.aspx?jobId=537) for information on Key Stage 2 marking costs.

Table 1: Descriptive Statistics

	Grade 2 (2009-2011)			Grade 5 (2009-2011)		
	Italy (1)	North/Centre (1)	South (1)	Italy (1)	North/Centre (1)	South (1)
A. Class Characteristics						
Female*	0.49 (0.50)	0.49 (0.50)	0.49 (0.50)	0.49 (0.50)	0.49 (0.50)	0.49 (0.50)
Immigrant*	0.10 (0.30)	0.14 (0.35)	0.03 (0.17)	0.10 (0.30)	0.14 (0.34)	0.03 (0.18)
Father HS*	0.34 (0.47)	0.34 (0.48)	0.33 (0.47)	0.32 (0.47)	0.33 (0.47)	0.30 (0.46)
Mother employed*	0.57 (0.49)	0.68 (0.47)	0.39 (0.49)	0.55 (0.50)	0.66 (0.47)	0.38 (0.49)
Pct correct: math	47.9 (14.6)	46.1 (12.9)	51.1 (16.7)	64.2 (12.9)	63.3 (10.9)	65.6 (15.5)
Pct correct: language	69.8 (10.9)	69.2 (9.2)	70.8 (13.3)	74.2 (8.9)	74.3 (7.5)	74.1 (10.8)
Class size	20.1 (3.40)	20.3 (3.35)	19.9 (3.48)	19.7 (3.72)	19.9 (3.67)	19.3 (3.76)
Score manipulation: math	0.06 (0.24)	0.02 (0.13)	0.14 (0.35)	0.07 (0.25)	0.02 (0.15)	0.14 (0.34)
Score manipulation: language	0.05 (0.23)	0.02 (0.14)	0.11 (0.31)	0.06 (0.23)	0.02 (0.15)	0.11 (0.31)
Number of classes	67,453	42,747	24,706	72,536	44,739	27,797
B. School Characteristics						
Number of classes	1.95 (1.10)	1.87 (1.01)	2.11 (1.27)	1.94 (1.10)	1.85 (0.98)	2.10 (1.28)
Enrollment	40.5 (25.2)	38.8 (23.0)	43.8 (28.6)	38.9 (25.2)	37.3 (22.8)	41.7 (28.9)
Number of schools	34,591	22,863	11,728	37,476	24,225	13,251
C. Institution Characteristics						
Number of schools	2.00 (1.05)	2.32 (1.13)	1.57 (0.74)	2.10 (1.09)	2.42 (1.17)	1.69 (0.81)
Number of classes	3.89 (1.97)	4.33 (1.95)	3.31 (1.85)	4.07 (1.95)	4.48 (1.91)	3.55 (1.88)
Enrollment	86.0 (40.6)	95.3 (39.5)	73.7 (38.7)	85.2 (40.5)	94.0 (39.1)	73.9 (39.3)
External monitor	0.22 (0.41)	0.20 (0.40)	0.23 (0.42)	0.22 (0.41)	0.20 (0.40)	0.23 (0.42)
Number of institutions	17,333	9,866	7,467	17,830	9,997	7,833

Notes: “Mean” and “s.d.” for class characteristics are computed using one observation per class; “Mean” and “s.d.” for school characteristics are computed using one observation per school; “Mean” and “s.d.” for institutions are computed using one observation per institution. \* conditional on non-missing survey response.

Table 2: OLS and IV/2SLS Estimates of the Effect of Class Size on Test Scores

	OLS			IV/2SLS					
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)	Italy (7)	North/Centre (8)	South (9)
A. Math									
Class size	-0.0078 (0.0070)	-0.0224*** (0.0067)	0.0091 (0.0146)	-0.0519*** (0.0134)	-0.0436*** (0.0115)	-0.0957*** (0.0362)	-0.0609*** (0.0196)	-0.0417** (0.0171)	-0.1294** (0.0507)
Enrollment	x	x	x	x	x	x	x	x	x
Enrollment squared	x	x	x	x	x	x	x	x	x
Interactions							x	x	x
N	140,010	87,498	52,512	140,010	87,498	52,512	140,010	87,498	52,512
B. Language									
Class size	0.0029 (0.0055)	-0.0188*** (0.0053)	0.0328*** (0.0114)	-0.0395*** (0.0106)	-0.0313*** (0.0092)	-0.0641** (0.0289)	-0.0409*** (0.0155)	-0.0215 (0.0136)	-0.0937** (0.0403)
Enrollment	x	x	x	x	x	x	x	x	x
Enrollment squared	x	x	x	x	x	x	x	x	x
Interactions							x	x	x
N	140,010	87,498	52,512	140,010	87,498	52,512	140,010	87,498	52,512

Notes: Columns 1-3 report OLS estimates of the effect of class size on scores. Columns 4-9 report 2SLS estimates using Maimonides' Rule as an instrument. The unit of observation is the class. Class size coefficients show the effect of 10 students. Models with interactions allow the quadratic running variable control to differ across windows of  $\pm 12$  students around each cutoff. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 3: OLS and IV/2SLS Estimates of the Effect of Class Size on Score Manipulation

	OLS			IV/2SLS					
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)	Italy (7)	North/Centre (8)	South (9)
A. Math									
Class size	-0.0163*** (0.0025)	-0.0074*** (0.0017)	-0.0309*** (0.0058)	-0.0186*** (0.0047)	-0.0042 (0.0031)	-0.0542*** (0.0143)	-0.0179*** (0.0069)	-0.0053 (0.0045)	-0.0471** (0.0202)
Enrollment	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>
Enrollment squared	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>
Interactions							<i>x</i>	<i>x</i>	<i>x</i>
N	139,996	87,491	52,505	139,996	87,491	52,505	139,996	87,491	52,505
B. Language									
Class size	-0.0166*** (0.0023)	-0.0120*** (0.0018)	-0.0244*** (0.0051)	-0.0202*** (0.0043)	-0.0116*** (0.0032)	-0.0400*** (0.0128)	-0.0161** (0.0063)	-0.0059 (0.0048)	-0.0379** (0.0177)
Enrollment	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>
Enrollment squared	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>	<i>x</i>
Interactions							<i>x</i>	<i>x</i>	<i>x</i>
N	140,003	87,493	52,510	140,003	87,493	52,510	140,003	87,493	52,510

Notes: Columns 1-3 report OLS estimates of the effect of class size on score manipulation. Columns 4-9 report 2SLS estimates using Maimonides' Rule as an instrument. Class size coefficients show the effect of 10 students. Models with interactions allow the quadratic running variable control to differ across windows of  $\pm 12$  students around each cutoff. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 4: Covariate Balance in the Monitoring Experiment

	Italy		North/Centre		South	
	Control Mean (1)	Treatment Difference (2)	Control Mean (3)	Treatment Difference (4)	Control Mean (5)	Treatment Difference (6)
A. Administrative Data on Schools						
Class size	19.812 [3.574]	0.0348 (0.0303)	20.031 [3.511]	0.0179 (0.0374)	19.456 [3.646]	0.0623 (0.0515)
Grade enrollment at school	53.119 [30.663]	-0.4011 (0.3289)	49.804 [27.562]	-0.5477 (0.3913)	58.483 [34.437]	-0.1410 (0.5909)
% in class sitting the test	0.939 [0.065]	0.0001 (0.0005)	0.934 [0.066]	0.0006 (0.0006)	0.947 [0.062]	-0.0007 (0.0008)
% in school sitting the test	0.938 [0.054]	-0.0001 (0.0005)	0.933 [0.055]	0.0005 (0.0006)	0.946 [0.051]	-0.0010 (0.0008)
% in institution sitting the test	0.937 [0.045]	-0.0001 (0.0004)	0.932 [0.043]	0.0005 (0.0005)	0.945 [0.045]	-0.0010 (0.0007)
B. Data Provided by School Staff						
Female students	0.482 [0.121]	0.0012 (0.0009)	0.483 [0.1179]	0.0004 (0.0011)	0.479 [0.126]	0.0027* (0.0016)
Immigrant students	0.097 [0.120]	0.0010 (0.0010)	0.137 [0.13]	0.0004 (0.0014)	0.031 [0.056]	0.0020*** (0.0007)
Father HS	0.25 [0.168]	0.0060*** (0.0016)	0.258 [0.163]	0.0061*** (0.0019)	0.238 [0.176]	0.0056** (0.0027)
Mother employed	0.441 [0.267]	0.0085*** (0.0024)	0.532 [0.258]	0.0067** (0.0031)	0.295 [0.210]	0.0117*** (0.0035)
C. Non-Response Indicators						
Missing data on father's education	0.223 [0.341]	-0.0217*** (0.0034)	0.225 [0.340]	-0.0186*** (0.0043)	0.221 [0.343]	-0.0271*** (0.0057)
Missing data on mother's occupation	0.195 [0.328]	-0.0168*** (0.0033)	0.196 [0.325]	-0.0083** (0.0042)	0.194 [0.333]	-0.0316*** (0.0054)
Missing data on country of origin	0.033 [0.163]	-0.0115*** (0.0013)	0.025 [0.143]	-0.0078*** (0.0014)	0.045 [0.192]	-0.0178*** (0.0026)
N	140,010		87,498		52,512	

Notes: Columns 1, 3 and 5 show means and standard deviations for variables listed at left. Other columns report coefficients from regressions of each variable on a treatment dummy (indicating classroom monitoring), grade and year dummies, and sampling strata controls (grade enrollment at institution, region dummies and their interactions). Standard deviations for the control group are in square brackets, robust standard errors are in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 5: Monitoring Effects on Score Manipulation and Test Scores

	Score manipulation			Test scores		
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)
A. Math						
Monitor at institution ( $M_{igkt}$ )	-0.029*** (0.002)	-0.010*** (0.001)	-0.062*** (0.004)	-0.112*** (0.006)	-0.075*** (0.005)	-0.180*** (0.012)
Means (sd)	0.064 (0.246)	0.020 (0.139)	0.139 (0.346)	0.007 (0.637)	-0.074 (0.502)	0.141 (0.796)
N	139,996	87,491	52,505	140,010	87,498	52,512
B. Language						
Monitor at institution ( $M_{igkt}$ )	-0.025*** (0.002)	-0.012*** (0.001)	-0.047*** (0.004)	-0.081*** (0.004)	-0.054*** (0.004)	-0.131*** (0.009)
Means (sd)	0.055 (0.229)	0.023 (0.149)	0.110 (0.313)	0.01 (0.523)	-0.005 (0.428)	0.035 (0.649)
N	140,003	87,493	52,510	140,010	87,498	52,512

Notes: Columns 1-3 report first stage estimates of the effect of a monitor at institution on score manipulation. Columns 4-6 show the reduced form effect of a monitor at institution on test scores. All models control for a quadratic in grade enrollment, segment dummies and their interactions. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values in these variables. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 6: IV/2SLS Estimates of the Effect of Class Size on Scores by Monitor at Institution

	Italy (1)	North/Centre (2)	South (3)
A. Math			
Class size* $M_{igkt}$	-0.0351 (0.0237)	-0.0389* (0.0211)	-0.0347 (0.0605)
Class size* (1- $M_{igkt}$ )	-0.0658*** (0.0207)	-0.0420** (0.0180)	-0.1433*** (0.0526)
$M_{igkt}$	-0.1736*** (0.0413)	-0.0815** (0.0376)	-0.3947*** (0.0959)
N	140,010	87,498	52,512
B. Language			
Class size* $M_{igkt}$	-0.0307 (0.0188)	-0.0208 (0.0169)	-0.0485 (0.0480)
Class size* (1- $M_{igkt}$ )	-0.0419** (0.0164)	-0.0212 (0.0144)	-0.0975** (0.0419)
$M_{igkt}$	-0.1033*** (0.0328)	-0.0545* (0.0300)	-0.2279*** (0.0764)
N	140,010	87,498	52,512

Notes: This table report 2SLS estimates using the interaction of Maimonides' Rule with monitor at institution ( $M_{igkt}$ ) as instruments. Class size coefficients show the effect of 10 students. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.



Table 7: Twin First Stages

	A. Score Manipulation					
	Math			Language		
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)
Maimonides' Rule ( $f_{igkt}$ )	-0.0009** (0.0004)	-0.0003 (0.0002)	-0.0019** (0.0009)	-0.0008** (0.0003)	-0.0003 (0.0003)	-0.0015** (0.0008)
Monitor at institution ( $M_{igkt}$ )	-0.029*** (0.002)	-0.010*** (0.001)	-0.062*** (0.004)	-0.025*** (0.002)	-0.012*** (0.001)	-0.047*** (0.004)
N	139,996	87,491	52,505	140,003	87,493	52,510
	B. Class size					
	Italy (1)	North/Centre (2)	South (3)			
Maimonides' Rule ( $f_{igkt}$ )	0.513*** (0.0006)	0.555*** (0.0008)	0.433*** (0.0011)			
Monitor at institution ( $M_{igkt}$ )	0.013 (0.024)	0.032 (0.027)	-0.009 (0.045)			
N	140,010	87,498	52,512			

Notes: Panel A report first stage estimates of the effect of the Maimonides' Rule and a monitor at institution on score manipulation. Panel B report first stage estimates of the effect of the Maimonides' Rule and a monitor at institution on class size. All models control for a quadratic in grade enrollment, segment dummies and their interactions. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values in these variables. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 8: IV/2SLS Estimates of the Effect of Class Size and Score Manipulation on Test Scores

	IV/2SLS			IV/2SLS (overidentified)			IV/2SLS (overidentified-interacted)		
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)	Italy (7)	North/Centre (8)	South (9)
A. Math									
Class size	0.0075 (0.0213)	-0.0029 (0.0298)	0.0062 (0.0441)	0.0024 (0.0190)	-0.0113 (0.0251)	0.0133 (0.0378)	0.0116 (0.0316)	0.0136 (0.0482)	0.0473 (0.0675)
Score manipulation	3.82*** (0.19)	7.33*** (0.79)	2.88*** (0.16)	3.82*** (0.19)	7.02*** (0.73)	2.87*** (0.16)	4.10*** (0.96)	9.21** (4.41)	3.33*** (0.86)
Class size * Score manipulation							-0.1464 (0.4814)	-1.2700 (2.1598)	-0.2273 (0.4304)
Overid test [P-value]				[0.914]	[0.600]	[0.541]	[0.914]	[0.475]	[0.476]
N	139,996	87,491	52,505	139,996	87,491	52,505	139,996	87,491	52,505
B. Language									
Class size	0.0121 (0.0173)	0.0049 (0.0196)	0.0127 (0.0385)	0.0218 (0.0153)	0.0109 (0.0174)	0.0491 (0.0329)	0.0325 (0.0308)	0.0098 (0.0320)	0.1337* (0.0800)
Score manipulation	3.29*** (0.18)	4.50*** (0.45)	2.80*** (0.18)	3.21*** (0.18)	4.34*** (0.42)	2.74*** (0.18)	3.59*** (1.03)	4.31* (2.25)	4.18*** (1.30)
Class size * Score manipulation							-0.2130 (0.4980)	-0.0029 (1.0898)	-0.7058 (0.6214)
Overid test [P-value]				[0.129]	[0.796]	[0.036]	[0.216]	[0.844]	[0.109]
N	140,003	87,493	52,510	140,003	87,493	52,510	140,003	87,493	52,510

Notes: Columns 1-3 show 2SLS estimates using Maimonides' Rule and monitor at institution as instruments. Columns 4-6 show overidentified 2SLS estimates which also use dummies for grade enrollment being in a 10 percent window below and above each cutoff (2 students) as instrument. Columns 7-9 add the interaction between class size and score manipulation and use the interaction of Maimonide's Rule with monitor at institution and the interactions of dummies for grade enrollment being in a 10 percent window below and above each cutoff with monitor at institution as instruments. Class size coefficients show the effect of 10 students. All models control for a quadratic in grade enrollment, segment dummies and their interactions. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values in these variables. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table 9: Maimonides' Rule and Covariate Balance

	Italy		North/Centre		South	
	Control Mean (1)	Treatment Difference (2)	Control Mean (3)	Treatment Difference (4)	Control Mean (5)	Treatment Difference (6)
A. Administrative Data on Schools						
% in class sitting the test	0.9392 [0.0643]	0.0000 (0.0001)	0.9345 [0.0657]	0.0001 (0.0001)	0.9471 [0.061]	0.0000 (0.0001)
% in school sitting the test	0.9386 [0.0534]	0.0001 (0.0001)	0.9339 [0.0548]	0.0001 (0.0001)	0.9464 [0.05]	0.0001 (0.0001)
% in institution sitting the test	0.9374 [0.0436]	-0.0001 (0.0001)	0.9327 [0.0426]	-0.0001 (0.0001)	0.9451 [0.0441]	-0.0000 (0.0001)
B. Data Provided by School Staff						
Female	0.482 [0.1205]	0.0000 (0.0002)	0.4836 [0.1176]	0.0002 (0.0002)	0.4792 [0.1251]	-0.0002 (0.0003)
Immigrant	0.0981 [0.1198]	-0.0007*** (0.0002)	0.1375 [0.1298]	-0.0007*** (0.0003)	0.0324 [0.0572]	-0.0004*** (0.0001)
Father HS	0.2546 [0.1678]	0.0006** (0.0003)	0.2613 [0.1626]	0.0002 (0.0003)	0.2434 [0.1755]	0.0013*** (0.0005)
Mother employed	0.4503 [0.2658]	0.0012*** (0.0004)	0.5356 [0.2574]	0.0010* (0.0005)	0.3082 [0.2138]	0.0016*** (0.0006)
C. Non-Response Indicators						
Missing data on father's education	0.2187 [0.3361]	0.0003 (0.0006)	0.2216 [0.3358]	0.0015** (0.0007)	0.2139 [0.3367]	-0.0018* (0.0010)
Missing data on mother's occupation	0.1925 [0.3239]	0.0002 (0.0006)	0.1963 [0.3231]	0.0014** (0.0007)	0.1861 [0.3251]	-0.0019* (0.0010)
Missing data on country of origin	0.0296 [0.1544]	-0.0001 (0.0002)	0.0232 [0.1361]	-0.0001 (0.0003)	0.0401 [0.1804]	-0.0000 (0.0005)
N	140,010		87,498		52,512	

Notes: Columns 1, 3 and 5 show means and standard deviations for variables listed at left. Other columns report coefficients from regressions of each variable on predicted class size (Maimonides' Rule), a quadratic in grade enrollment, segment dummies and their interactions, grade and year dummies, and sampling strata controls (grade enrollment at institution, region dummies and their interactions). Standard deviations for the control group are in square brackets, robust standard errors are in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Figure 1: Manipulation Rates by Province

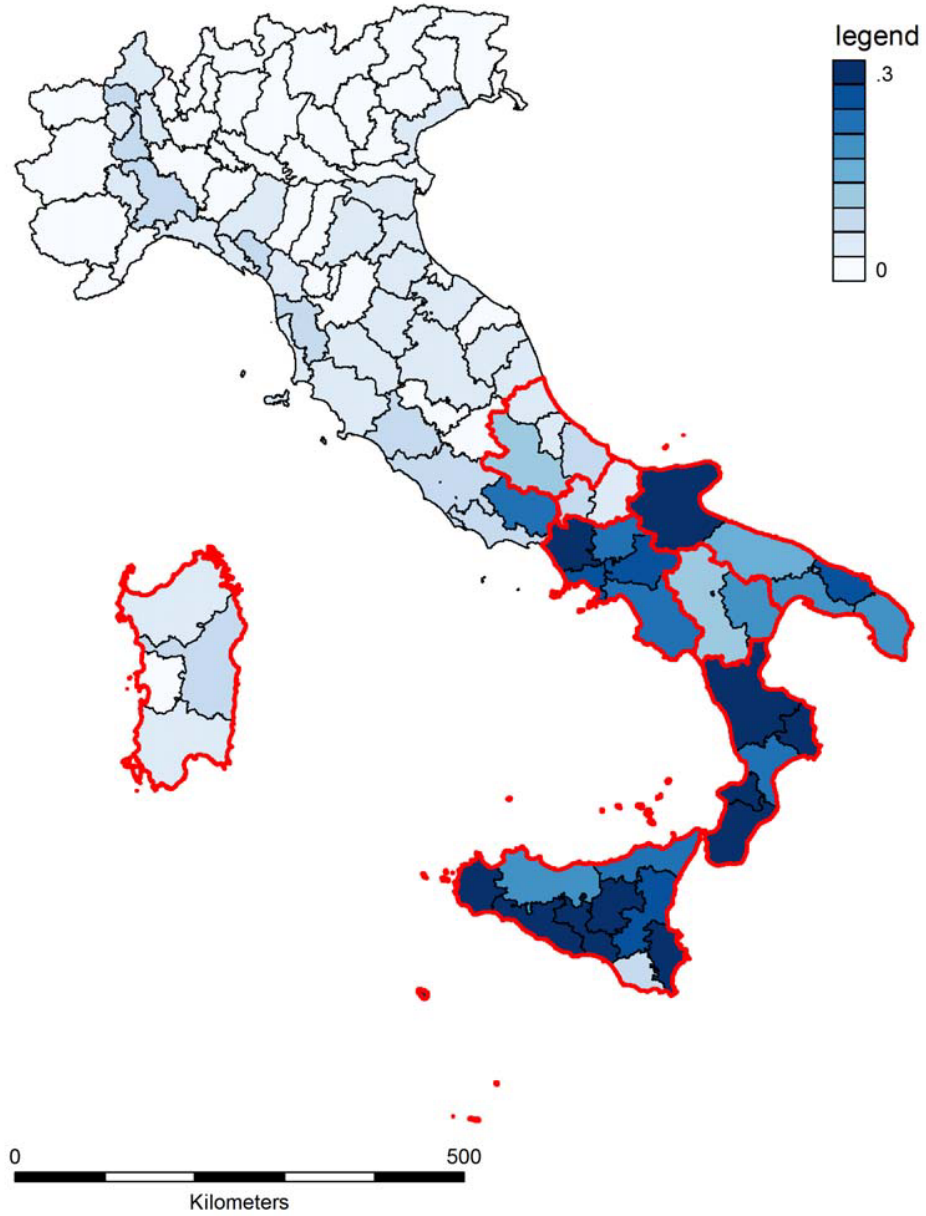
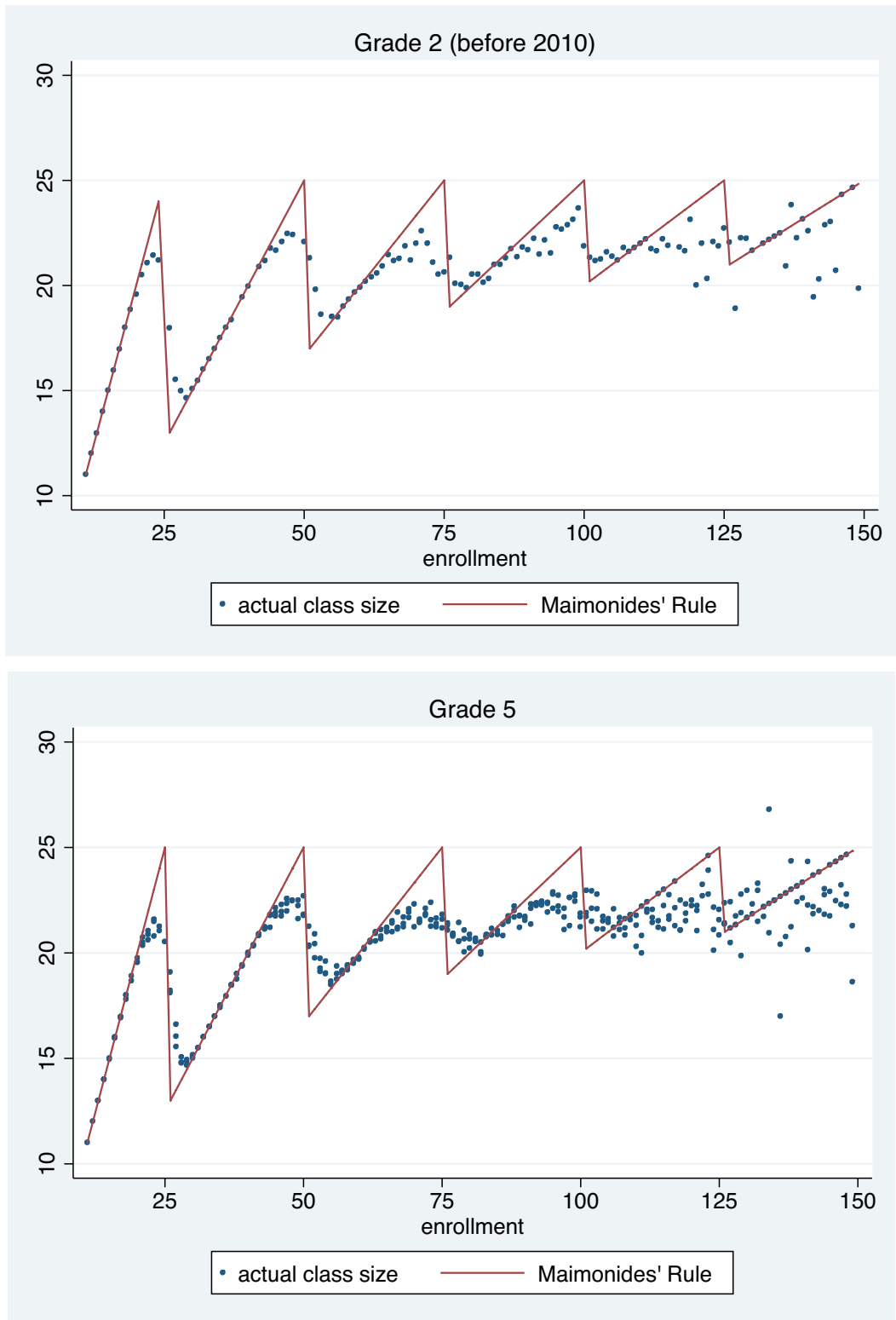


Figure 2: Class Size by Enrollment in Pre-reform Years



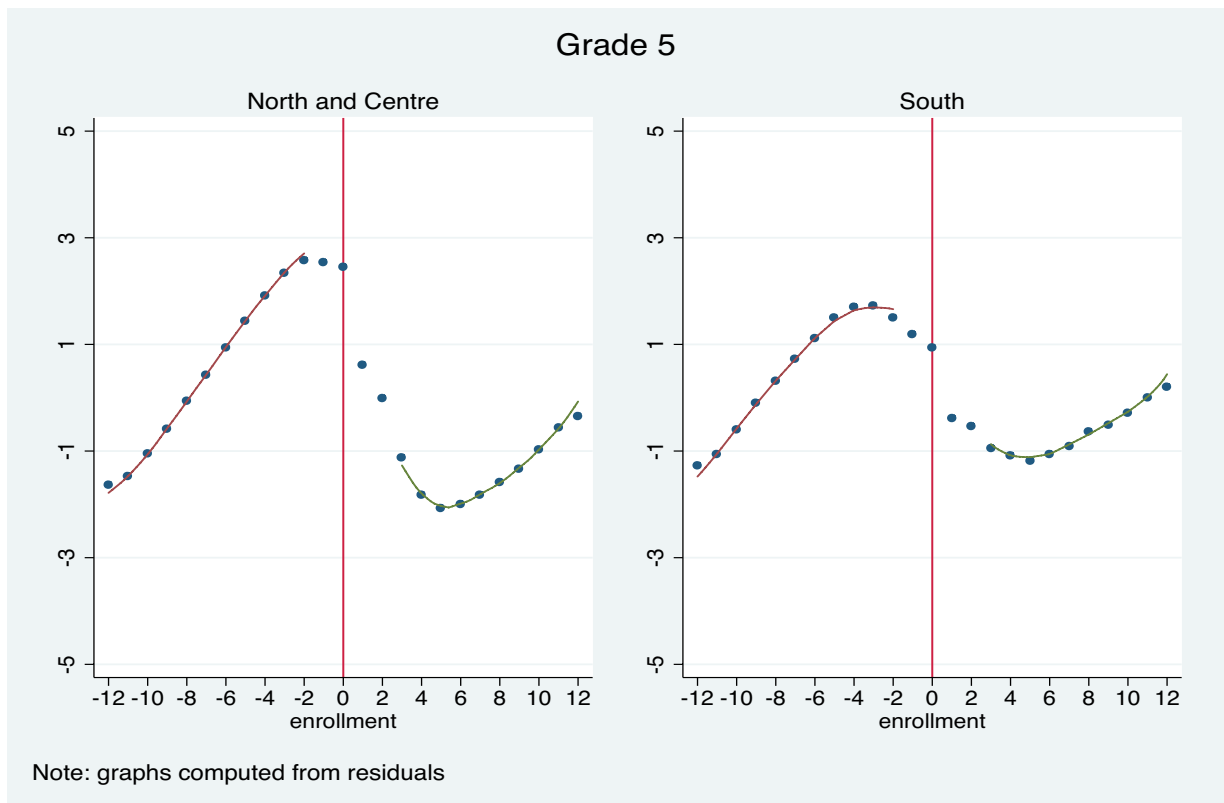
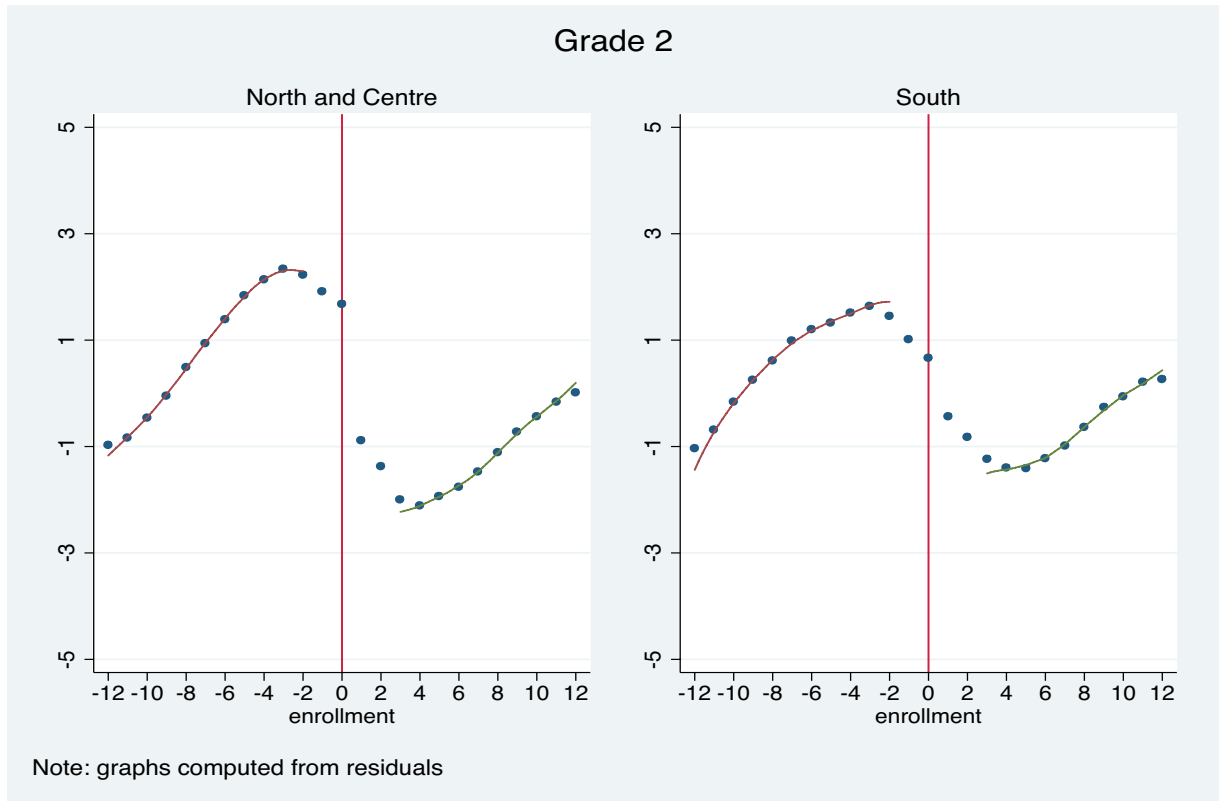
Notes: The figure shows actual class size and as predicted by Maimonides' Rule in pre-reform years

Figure 3: Class Size by Enrollment in Post-reform Years



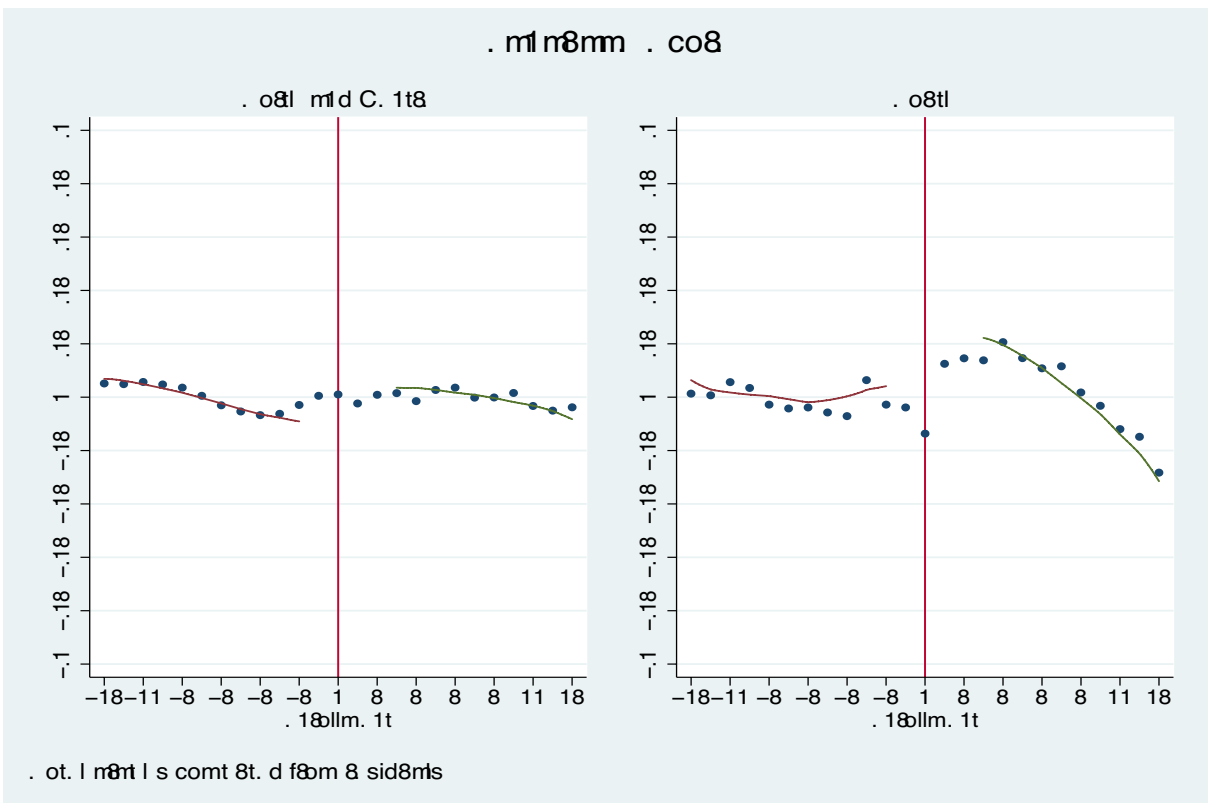
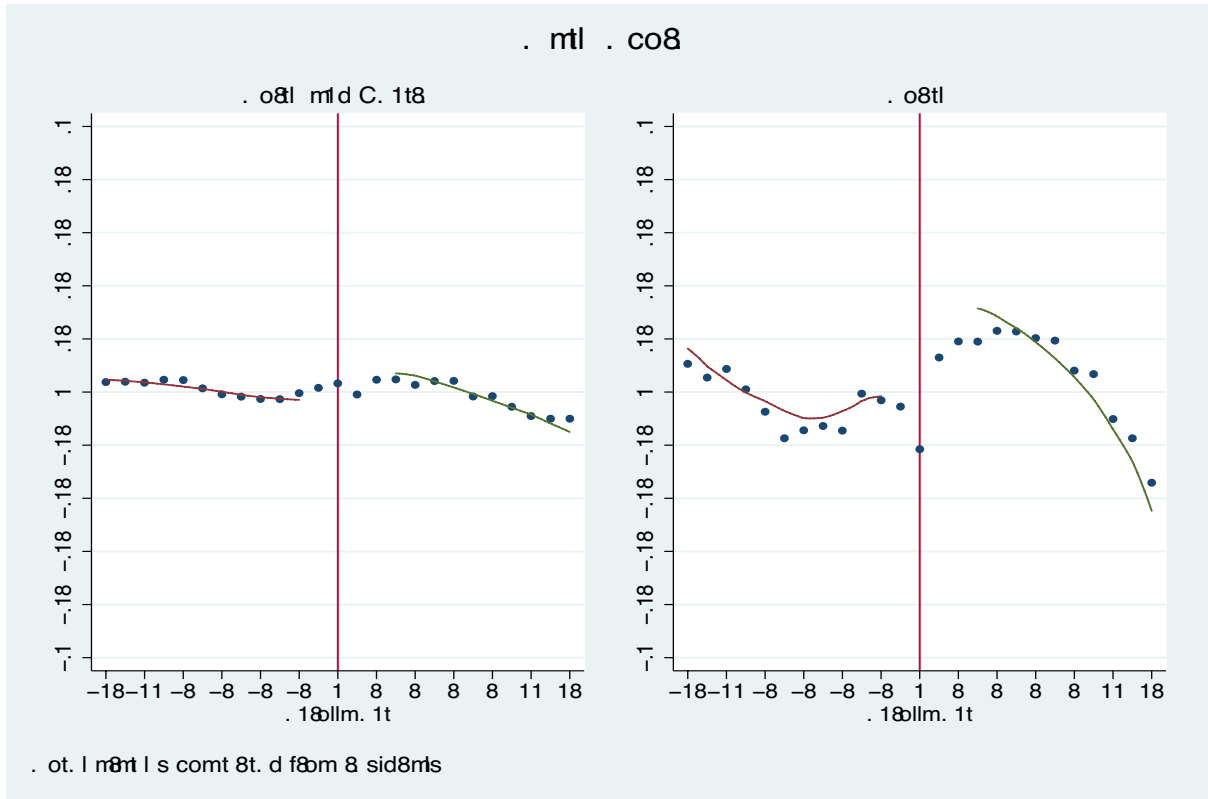
Notes: The figure shows actual class size and as predicted by Maimonides' Rule in post-reform years

Figure 4: Class Size and Enrollment, centered at Maimonides Cutoffs



Notes: The solid line shows a one-sided LLR fit.

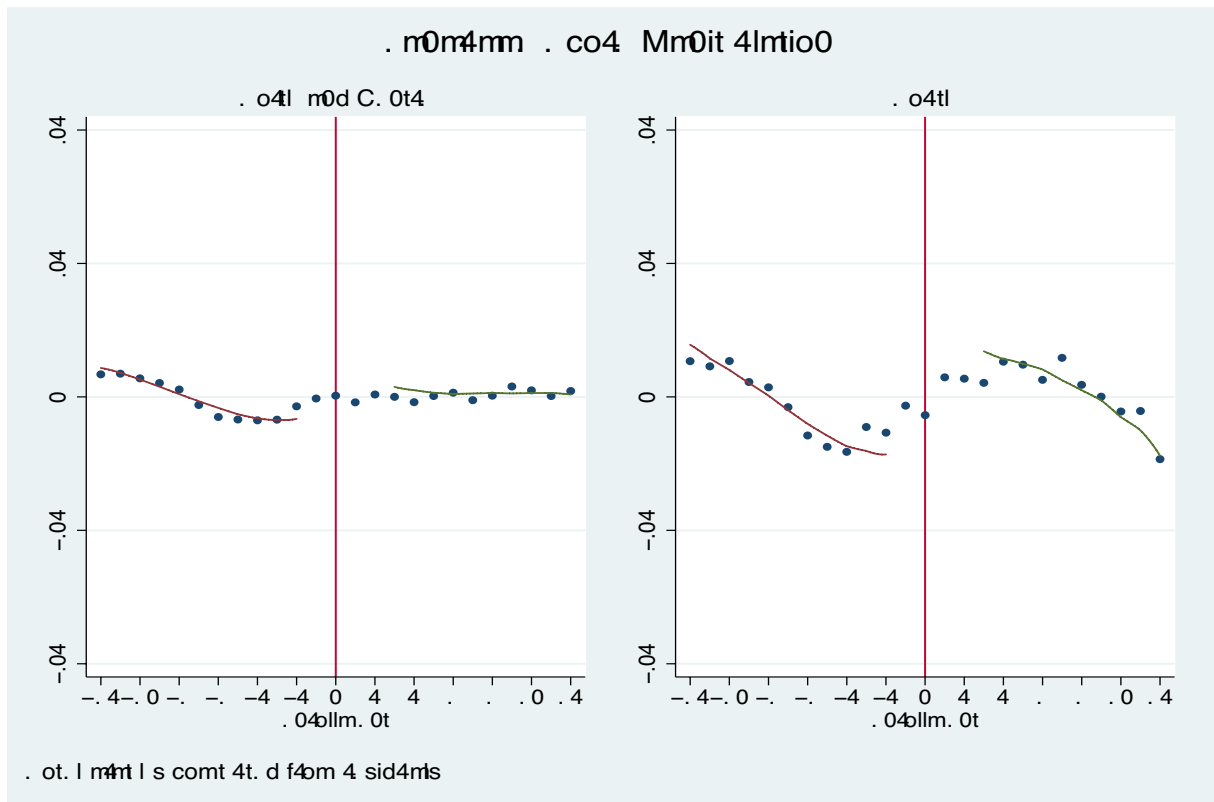
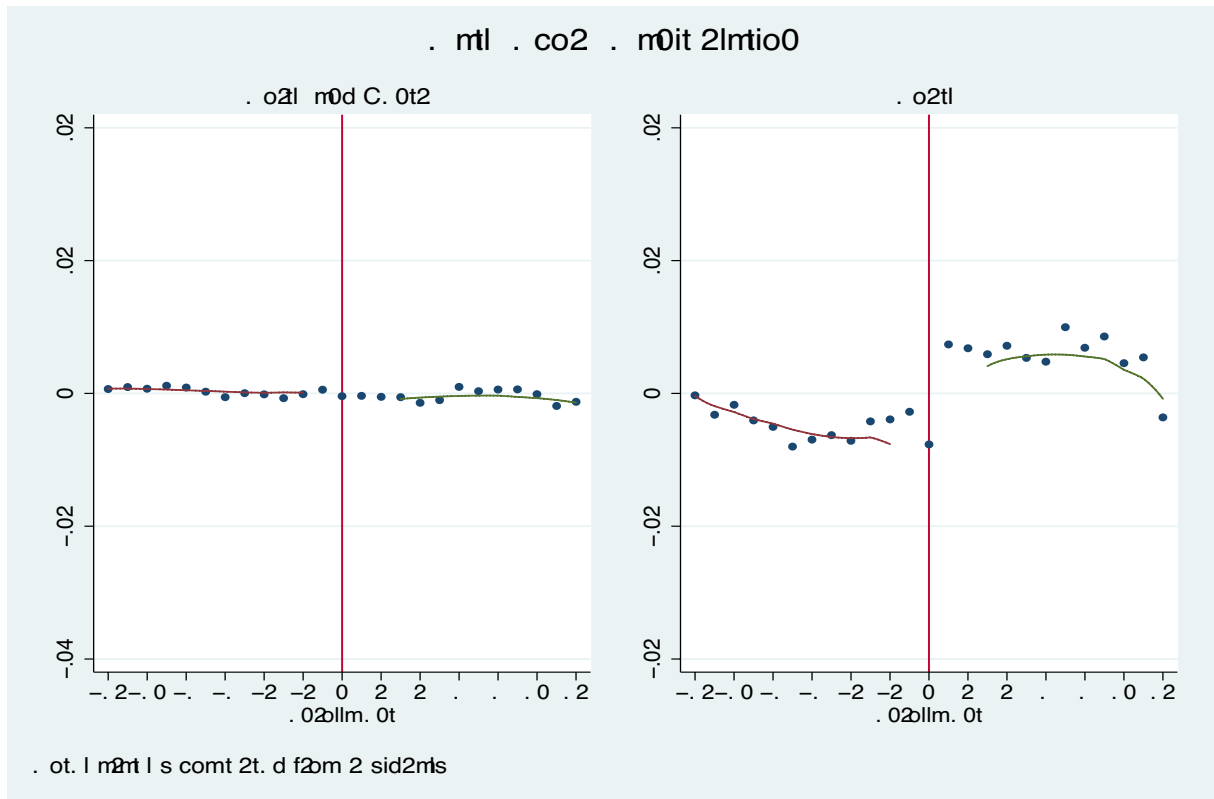
Figure 5: Test Scores and Enrollment, centered at Maimonides Cutoffs



Notes: The solid line shows a one-sided LLR fit.

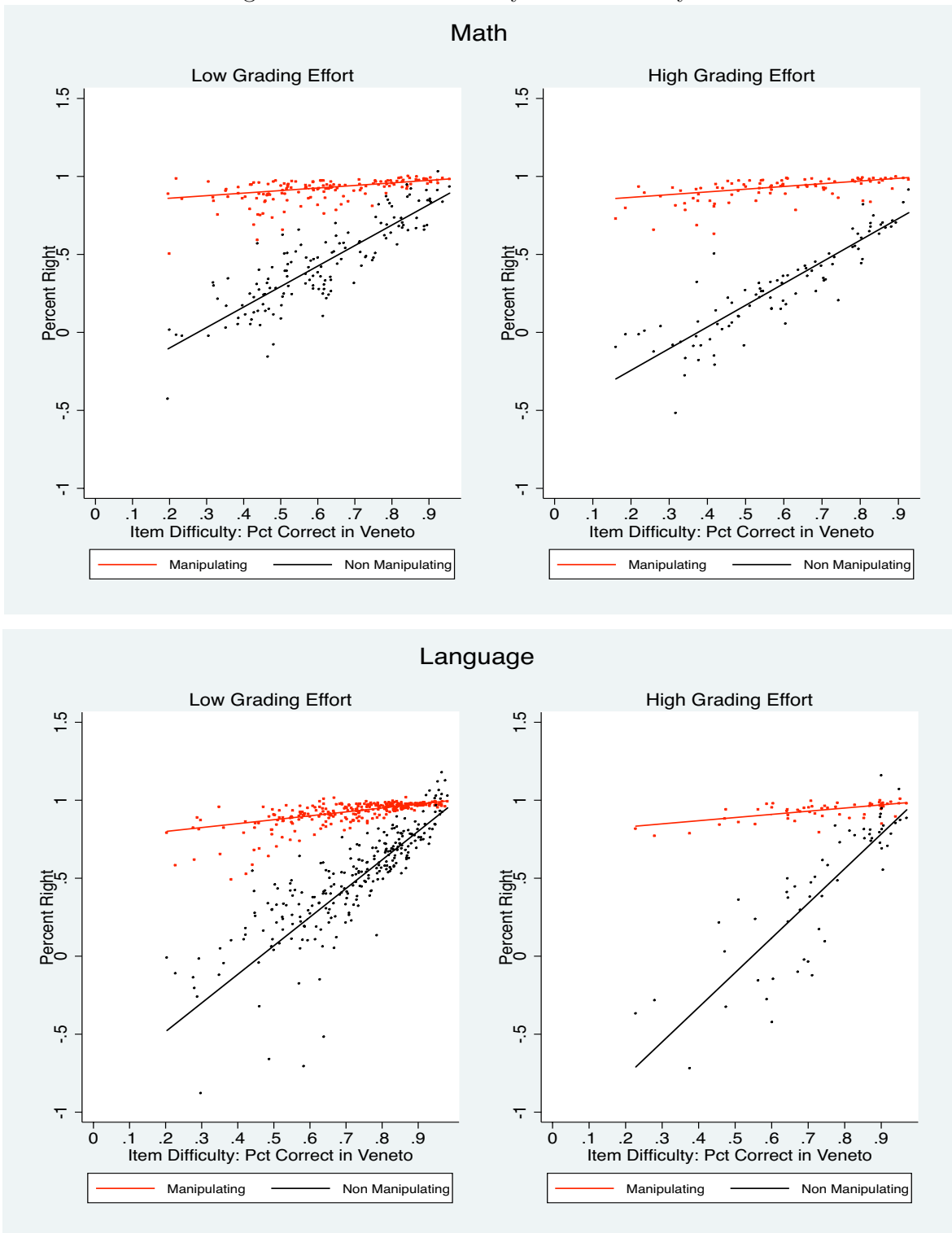


Figure 6: Score Manipulation and Enrollment, centered at Maimonides Cutoffs



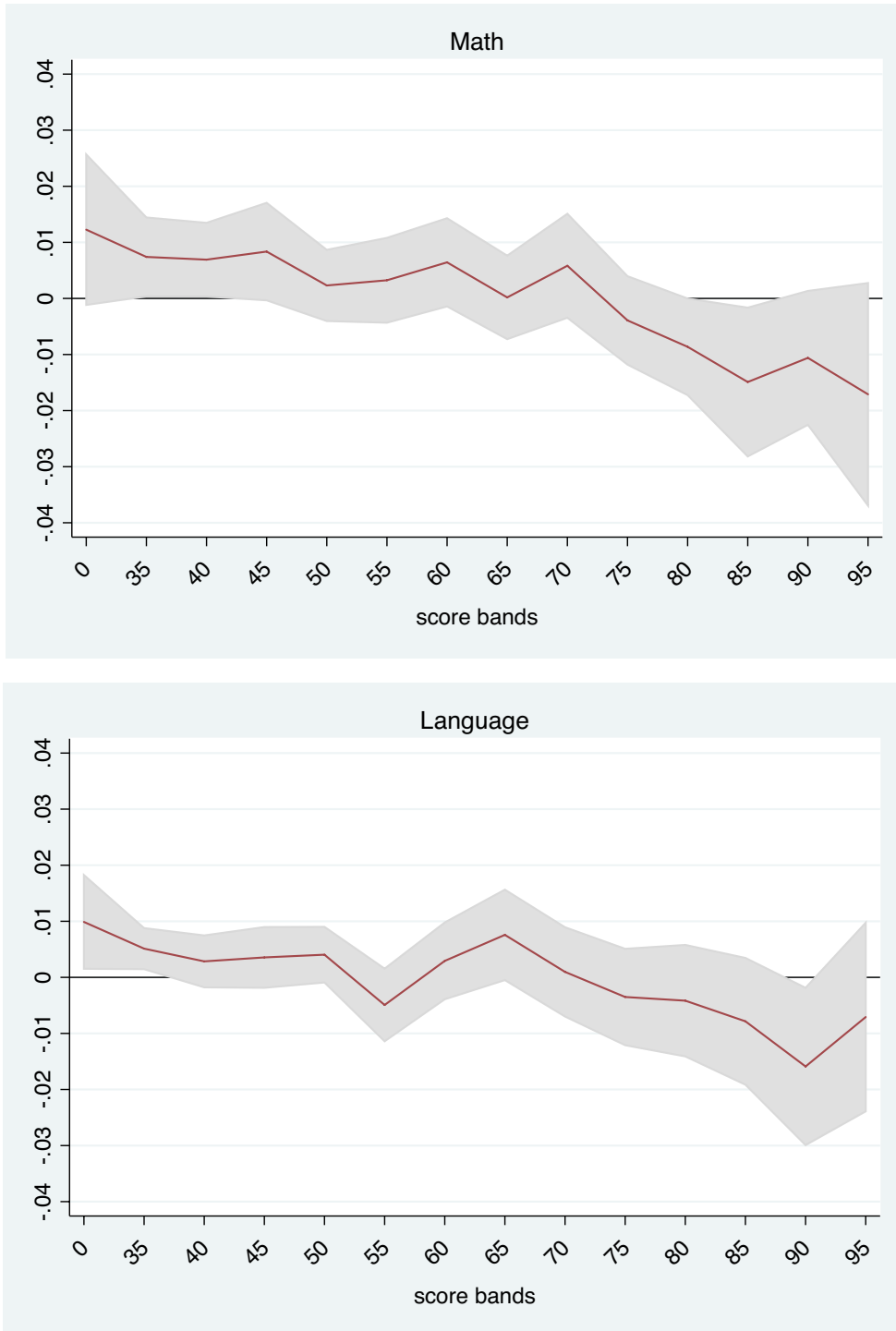
Notes: The solid line shows a one-sided LLR fit.

Figure 7: Score Gradient by Item Difficulty



Notes: The figures plot the average potential score on item  $j$  under manipulation for complying classes and the average potential score on item  $j$  without manipulation for the same classes against the percent correct answers in monitored institutions in Veneto. The sample is restricted to the South.

Figure 8: Class Size Effects on Score Distributions



Notes: These figures plot 2SLS estimates of class size effects on indicators for 5-point score bands using data from the South. Models and methods are the same as those used to construct the estimated class size effects in Table 2

## References

- ABADIE, A. (2002): “Bootstrap tests for distributional treatment effects in instrumental variables models,” *Journal of the American Statistical Association*, 97, 284–292.
- ANGRIST, J. D., AND V. LAVY (1999): “Using Maimonides’ Rule to Estimate the Effect of Class Size on Scholastic Achievement,” *Quarterly Journal of Economics*, 114(2), 533–575.
- ANGRIST, J. D., P. PATHAK, AND C. R. WALTERS (2013): “Explaining Charter School Effectiveness,” *American Economic Journal: Applied Economics*, 5(4), 1–27.
- BAKER, O., AND D. PASERMAN (2013): “Grade Enrollment Sorting under an Incentives-Based Class Size Reduction Program,” Unpublished mimeo.
- BALLATORE, R., M. FORT, AND A. ICHINO (2013): “The Tower of Babel in the classroom: immigrants and natives in Italian schools,” Unpublished mimeo.
- BANERJEE, A., AND E. DUFLO (2006): “Addressing Absence,” *Journal of Economic Perspectives*, 20(1), 117–132.
- BERTONI, M., G. BRUNELLO, AND L. ROCCO (2013): “When the cat is near, the mice won’t play: The effect of external examiners in Italian schools,” *Journal of Public Economics*, 104, 65–77.
- BEZDEK, J. (1981): *Pattern Recognition with Fuzzy Objective Function Algorithms*. Plenum Press, New York.
- BÖHLMARK, A., AND M. LINDAHL (2013): “Independent Schools and Long-Run Educational Outcomes - Evidence from Sweden’s Large Scale Voucher Reform,” forthcoming, *Economica*.
- BONESRONNING, H. (2003): “Class size effects on student achievement in Norway: Patterns and explanations,” *Southern Economic Journal*.
- BRATTI, M., D. CHECCHI, AND A. FILIPPIN (2007): “Territorial differences in Italian students’ mathematical competences: Evidence from PISA,” *Giornale degli Economisti e Annali di Economia*, 66(3), 299–335.

- BRUNELLO, G., AND D. CHECCHI (2005): “School quality and family background in Italy,” *Economics of Education Review*, 24, 563–577.
- CARD, D., D. S. LEE, Z. PEI, AND A. WEBER (2012): “Nonlinear Policy Rules and the Identification and Estimation of Causal Effects in a Generalized Regression Kink Design,” NBER Working Paper, 18564.
- CHAUDHURY, N., J. HAMMER, M. KREMER, K. MURALIDHARAN, AND F. H. ROGERS (2006): “Missing in Action: Teacher and Health Worker Absence in Developing Countries,” *Journal of Economic Perspectives*, 20(1), 91–116.
- CHETTY, R., J. FRIEDMAN, N. HILGER, E. SAEZ, D. SCHANZENBACH, AND D. YAGAN (2011): “How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project STAR,” *Quarterly Journal of Economics*, 126(4), 1593–1660.
- CLOTFELTER, C. T., H. F. LADD, AND J. L. VIGDOR (2009): “Are Teacher Absences Worth Worrying About in the United States?,” *Education Finance and Policy*, 4(2), 115–149.
- COSTANTINI, M., AND C. LUPI (2006): “Divergence and long-run equilibria in Italian regional unemployment,” *Applied Economics Letters*, 13(14), 899–904.
- DEE, T. S., B. A. JACOB, J. MCCRARY, AND J. ROCKOFF (2011): “Rules and Discretion in the Evaluation of Students and Schools: The Case of the New York Regents Examinations,” Columbia Business School Research Paper. Available at SSRN: <http://ssrn.com/abstract=1915387>.
- DEPAOLA, M., V. SCOPPA, AND V. PUPO (2014): “Absenteeism in the Italian Public Sector: The Effects of Changes in Sick Leave Policy,” *Journal of Labor Economics*, 32(2), 337–360.
- DOBBELSTEEN, S., J. LEVIN, AND H. OOSTERBEEK (2002): “The causal effect of class size on scholastic achievement: Distinguishing the pure class size effect from the effect of changes in class composition,” *Oxford Bulletin of Economics and Statistics*, 64(1), 17–38.
- GARY-BOBO, R. J., AND M.-B. MAHJOUR (2006): “Estimation of class-size effects, using Maimonides’ rule: the case of French junior high schools,” CEPR Discussion Papers 5754.

- GUIO, L., P. SAPIENZA, AND L. ZINGALES (2004): “The Role of Social Capital in Financial Development,” *American Economic Review*, 94(3), 526–556.
- (2010): “Civic Capital as the Missing Link,” in *Handbook of Social Economics*, ed. by A. B. Jess Benhabib, and M. Jackson. North Holland.
- HANUSHEK, E. A. (1995): “Interpreting recent research on schooling in developing countries,” *The World Bank Research Observer*, X, 227–246.
- HOLMSTROM, B., AND P. MILGROM (1991): “Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design,” *Journal of Law, Economics, and Organization*, 7, 24–52.
- HOXBY, C. (2000): “Peer Effects in the Classroom: Learning from Gender and Race Variation,” NBER Working paper 7867.
- ICHINO, A., AND P. ICHINO (1997): “Culture, Discrimination and Individual Productivity: Regional Evidence from Personnel Data in a Large Italian Firm,” CEPR Discussion Papers 1709.
- ICHINO, A., AND G. MAGGI (2000): “Work Environment and Individual Background: Explaining Regional Shirking Differentials in a Large Italian Firm,” *Quarterly Journal of Economics*, 115(3), 933–959.
- ICHINO, A., AND G. TABELLINI (2014): “Freeing the Italian School System,” forthcoming, *Labour Economics*.
- IMBENS, G., AND K. KALYANARAMAN (2012): “Optimal Bandwidth Choice for the Regression Discontinuity Estimator,” *Review of Economic Studies*, 79(3), 933–959.
- INVALSI (2010): “Sistema Nazionale di Valutazione - A.S. 2009/2010, Rilevazione degli apprendimenti,” *Technical Report*.
- JACOB, B., AND S. LEVITT (2003): “Rotten Apples: An Investigation of the Prevalence and Predictors of Teacher Cheating,” *Quarterly Journal of Economics*, 118(3), 843–77.

- KANE, T. J., C. E. ROUSE, AND D. STAIGER (1999): “Estimating Returns to Schooling When Schooling is Misreported,” NBER Working Paper 7235.
- KRUEGER, A. (1999): “Experimental estimates of education production functions,” *Quarterly Journal of Economics*, 114, 497–532.
- LEUVEN, E., H. OOSTERBEEK, AND M. RONNING (2008): “Quasi-experimental estimates of the effect of class size achievement in Norway,” *The Scandinavian Journal of Economics*, 110(4), 663–693.
- LEWBEL, A. (2007): “Estimation of Average Treatment Effects with Misclassification,” *Econometrica*, 2(3), 537–551.
- MAHAJAN, A. (2006): “Identification and Estimation of Regression Models with Misclassification,” *Econometrica*, 74(3), 631–665.
- NANNICINI, T., A. STELLA, G. TABELLINI, AND U. TROIANO (2013): “Social Capital and Political Accountability,” *American Economic Journal: Economic Policy*, 5, 1957–1969.
- NEAL, D. (2013): “The Consequences of Using One Assessment System to Pursue Two Objectives,” NBER Working paper 19214.
- PIKETTY, T. (2004): “Should we reduce class size or school segregation? Theory and evidence from France,” presentation at the Roy Seminars, Association pour le développement de la recherche en économie et en statistique (ADRES), 22 November, available at: <http://www.adres.polytechnique.fr/SEMINAIRE/221104b.pdf>.
- PUTNAM, R., R. LEONARDI, AND R. NANETTI (1993): *Making Democracy Work*. Princeton University Press, Princeton.
- QUINTANO, C., R. CASTELLANO, AND S. LONGOBARDI (2009): “A Fuzzy Clustering Approach to Improve the Accuracy of Italian Student Data. An Experimental Procedure to Correct the Impact of the Outliers on Assessment Test Scores,” *Statistica & Applicazioni*, Vol.VII(2), 149–171.

SEVERSON, K. (2011): “Systematic Cheating Is Found in Atlanta’s School System,” *New York Times*, July 11, Accessed at: <http://www.nytimes.com/2011/07/06/education/06atlanta.html>.

URQUIOLA, M., AND E. VERHOOGEN (2009): “Class size caps, sorting, and the regression discontinuity design,” *American Economic Review*, 99(1), 179–215.

WOESSMANN, L. (2005): “Educational production in Europe,” *Economic Policy*, 43, 445–493.



## Appendix

### Score Manipulation Imputation

Our imputation is closely related to that used by INVALSI and described in Quintano et al. (2009). INVALSI assigns a manipulation probability to each class in three steps.

The first step computes the following four summary statistics.

(1) Within-class average score:

$$\bar{p}_i = \frac{\sum_{j=1}^{N_i} p_{ji}}{N_i}, \quad (10)$$

where  $p_{ji}$  denotes the score of student  $j$  in class  $i$ ;  $N_i$  denotes the number of test-takers in class  $i$ .

(2) Within-class standard deviation of scores:

$$\sigma_i = \sqrt{\frac{\sum_{j=1}^{N_i} (p_{ji} - \bar{p}_i)^2}{N_i}}. \quad (11)$$

(3) Within-class average percent missing

$$MC_i = \frac{\sum_{j=1}^{N_i} M_{ji}}{N_i}, \quad (12)$$

where  $M_{ji}$  is the fraction of test items skipped by student  $j$  in class  $i$ .

(4) Within-class index of answer homogeneity:

$$\bar{E}_i = \frac{\sum_{q=1}^Q E_{qi}}{Q}, \quad (13)$$

where  $q = 1, \dots, Q$  indexes test items and  $E_{qi}$  is a Gini measure of homogeneity that equals value zero if all students in class  $i$  provide the same answer to item  $q$ . This can be interpreted as the Herfindahl index of the share of students with similar response patterns in the class.

In the second step, the first two principal components are extracted from the  $4 \times 4$  correlation matrix determined by these indicators, yielding a percentage of explained variance which is - across years, subjects and grades - well above 90%. Denote these principal com-

ponents by  $\psi_{1i}$  and  $\psi_{2i}$ . The third step consists of a cluster analysis that creates  $G$  groups from the distribution of  $(\psi_{1i}, \psi_{2i})$ . INVALSI sets  $G = 8$ , yielding a matrix whose elements are, for each class, eight group membership probabilities. This procedure is known as “fuzzy clustering” (see Bezdek, 1981), since data elements (classes, in our setting) can be assigned to one or more groups. With “hard clustering”, data elements belong to exactly one cluster.

INVALSI identifies likely manipulators as those in the group with values of  $(\psi_{1i}, \psi_{2i})$  that are most extreme (see Figure 8 in Quintano et al. 2009). In practice, the suspicious group is characterized by (i) abnormally large values of  $\bar{p}_i$ , and (ii) small values of  $\sigma_i$ ,  $MC_i$  and  $\bar{E}_i$ , relative to the population average of these indicators. This group is flagged as the “outlier” or manipulating cluster. The INVALSI manipulation indicator gives, for each class, the membership probability for this cluster. Our hard clustering computations codes a dummy for manipulating classes. This dummy indicates classes whose values of  $(\psi_{1i}, \psi_{2i})$  belong to the manipulating cluster identified by INVALSI.

Table A1: Reduced Form Estimates of the Effect of Maimonides' Rule on Class Size, Test Scores, and Score Manipulation

	Math			Language		
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)
A. Class size						
Maimonides' Rule	0.513*** (0.006)	0.555*** (0.008)	0.433*** (0.011)			
Means (sd)	19.88 (3.58)	20.07 (3.52)	19.58 (3.64)			
N	140,010	87,498	52,512			
B. Test Scores						
Maimonides' Rule	-0.0031*** (0.0010)	-0.0023** (0.0009)	-0.0056** (0.0022)	-0.0021*** (0.0008)	-0.0012 (0.0008)	-0.0041** (0.0017)
Means (sd)	0.007 (0.637)	-0.074 (0.502)	0.141 (0.796)	0.01 (0.523)	-0.005 (0.428)	0.035 (0.649)
N	140,010	87,498	52,512	140,010	87,498	52,512
C. Score Manipulation						
Maimonides' Rule	-0.0009*** (0.0004)	-0.0003 (0.0002)	-0.0020** (0.0009)	-0.0008** (0.0003)	-0.0003 (0.0003)	-0.0016** (0.0008)
Means (sd)	0.065 (0.246)	0.02 (0.139)	0.139 (0.346)	0.055 (0.229)	0.023 (0.149)	0.110 (0.313)
N	139,996	87,491	52,505	140,003	87,493	52,510

Notes: This table shows the reduced form effect of the Maimonides' Rule on class size (Panel A), test scores (Panel B), score manipulation (Panel C). All models control for a quadratic in grade enrollment, segment dummies and their interactions. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF grade and year dummies, and dummies for missing values in these variables. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Table A2: First Stage Estimates for Over-Identified Models

	Class size			Score manipulation math			Score manipulation language		
	Italy	North/Centre	South	Italy	North/Centre	South	Italy	North/Centre	South
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Maimonides' Rule ( $f_{igkt}$ )	0.704*** (0.0059)	0.753*** (0.0069)	0.617*** (0.0107)	-0.0009** (0.0005)	-0.0003 (0.0003)	-0.0021* (0.0011)	-0.0014*** (0.0004)	-0.0008** (0.0003)	-0.0024** (0.0010)
Monitor at institution ( $M_{igkt}$ )	0.010 (0.023)	0.029 (0.026)	-0.013 (0.044)	-0.029*** (0.002)	-0.010*** (0.001)	-0.062*** (0.004)	-0.025*** (0.002)	-0.012*** (0.001)	-0.047*** (0.004)
2 students below cutoff	-1.427*** (0.083)	-1.154*** (0.101)	-1.865*** (0.138)	0.002 (0.005)	-0.002 (0.003)	0.008 (0.012)	0.010** (0.005)	0.005 (0.004)	0.018 (0.011)
1 student below cutoff	-2.258*** (0.093)	-2.053*** (0.116)	-2.580*** (0.150)	0.001 (0.005)	0.001 (0.004)	0.000 (0.012)	0.007 (0.005)	0.009** (0.004)	0.002 (0.011)
1 student above cutoff	2.411*** (0.097)	3.026*** (0.132)	1.519*** (0.138)	0.000 (0.006)	0.003 (0.005)	-0.004 (0.013)	-0.001 (0.005)	-0.001 (0.004)	-0.001 (0.012)
2 students above cutoff	1.247*** (0.083)	1.546*** (0.114)	0.826*** (0.120)	0.001 (0.006)	-0.004 (0.004)	0.007 (0.013)	-0.007 (0.005)	-0.005 (0.004)	-0.012 (0.009)
N	140,010	87,498	52,512	139,996	87,491	52,505	140,003	87,493	52,510

Notes: Columns 1-3 report first stage estimates of the effect of the Maimonides' Rule, a monitor at institution and dummies for grade enrollment being in a 10 percent window below and above each cutoff on class size. Columns 4-9 show first stage estimates of the effect of the Maimonides' Rule, a monitor at institution and dummies for grade enrollment being in a 10 percent window (2 students) above and below each cutoff on score manipulation. All models control for a quadratic in grade enrollment, segment dummies and their interactions. The unit of observation is the class. Robust standard errors, clustered on school and grade, are shown in parentheses. Control variables include: % female students, % immigrants, % fathers at least high school graduate, % employed mothers, % unemployed mothers, % mother NILF, grade and year dummies, and dummies for missing values in these variables. All regressions include sampling strata controls (grade enrollment at institution, region dummies and their interactions). \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

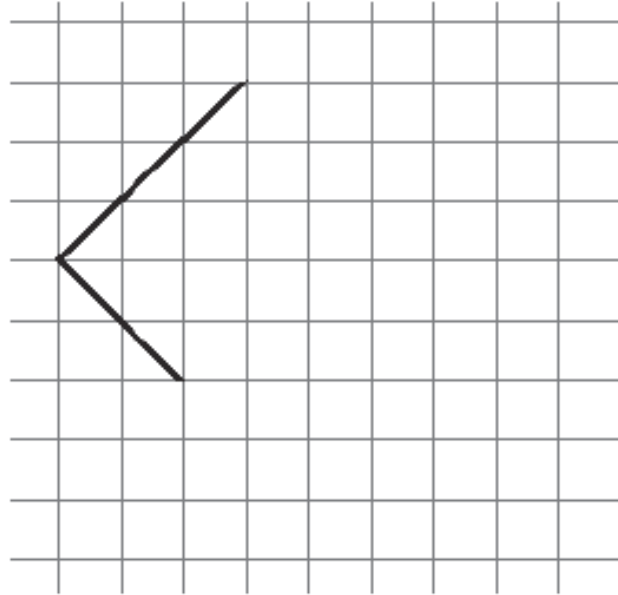
Table A3: Covariates and Maimonides' Rule with and without External Monitors

	Institutions with Monitor			Institutions without Monitor		
	Italy (1)	North/Centre (2)	South (3)	Italy (4)	North/Centre (5)	South (6)
	A. Administrative Data on Schools					
% in class sitting the test	0.0001 (0.0002)	0.0002 (0.0002)	0.0000 (0.0003)	0.0000 (0.0001)	0.0000 (0.0001)	0.0000 (0.0002)
% in school sitting the test	0.0003 (0.0002)	0.0003 (0.0002)	0.0002 (0.0003)	0.0001 (0.0001)	0.0001 (0.0001)	0.0001 (0.0002)
% in institution sitting the test	-0.0000 (0.0001)	-0.0000 (0.0002)	0.0001 (0.0003)	-0.0001* (0.0001)	-0.0002* (0.0001)	-0.0000 (0.0001)
	B. Data Provided by School Staff					
Female	-0.0003 (0.0003)	-0.0006 (0.0004)	0.0001 (0.0006)	0.0001 (0.0002)	0.0005* (0.0002)	-0.0003 (0.0003)
Immigrant	-0.0005 (0.0003)	-0.0002 (0.0005)	-0.0007** (0.0003)	-0.0007*** (0.0002)	-0.0009*** (0.0003)	-0.0003* (0.0002)
Father HS	-0.0005 (0.0005)	-0.0002 (0.0006)	-0.0014 (0.0010)	0.0010*** (0.0003)	0.0003 (0.0004)	0.0020*** (0.0005)
Mother employed	0.0001 (0.0008)	0.0003 (0.0010)	-0.0004 (0.0012)	0.0015*** (0.0004)	0.0012** (0.0006)	0.0022*** (0.0006)
	C. Non-Response Indicators					
Missing data on father's education	0.0014 (0.0011)	0.0012 (0.0013)	0.0019 (0.0020)	0.0000 (0.0007)	0.0016** (0.0008)	-0.0026** (0.0012)
Missing data on mother's occupation	0.0018* (0.0011)	0.0017 (0.0013)	0.0020 (0.0019)	-0.0002 (0.0007)	0.0012 (0.0008)	-0.0028** (0.0011)
Missing data on country of origin	0.0006 (0.0004)	0.0003 (0.0004)	0.0011 (0.0008)	-0.0002 (0.0003)	-0.0002 (0.0003)	-0.0003 (0.0006)
N	34,325	22,174	12,151	105,685	65,324	40,361

Notes: This table reports coefficients from regressions of the variables listed at left on Maimonides' Rule, controlling for a quadratic in grade enrollment, enrollment segment dummies and their interactions, grade and year dummies, and sampling strata controls (grade enrollment at institution, region dummies and their interactions). Columns 1-3 show results for the sample with monitors; columns 4-6 show results for the sample without monitors. Robust standard errors, clustered on school and grade, are shown in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

Figure A1: Example of open-ended question in math test - V grade 2010/11

**D23. Osserva la seguente figura.**



- a. Completa la figura in modo da ottenere un quadrato.**
- b. Spiega come hai fatto per disegnare il quadrato.**

.....

.....

.....

Figure A2: Example of open-ended question in language test - V grade 2010/11

**C4. Nella frase che segue inserisci le parole mancanti scegliendole da questa lista: *così, dove, perché, però, se, siccome*.**

..... non conoscevo la strada, ho chiesto a una signora .....  
dovevo andare; ..... non mi sono perso.





**This working paper has been produced by  
the School of Economics and Finance at  
Queen Mary University of London**

**Copyright © 2015 Joshua D. Angrist, Erich Battistin  
and Daniela Vuri. All rights reserved**

**School of Economics and Finance  
Queen Mary University of London  
Mile End Road  
London E1 4NS  
Tel: +44 (0)20 7882 7356  
Fax: +44 (0)20 8983 3580  
Web: [www.econ.qmul.ac.uk/research/workingpapers/](http://www.econ.qmul.ac.uk/research/workingpapers/)**